

Preliminary results on Pulmonary Tuberculosis detection in Chest X-Ray using Convolutional Neural Networks

Márcio Eloi Colombo Filho¹[0000-0003-3779-0192] ,
Rafael Mello Galliez²[0000-0003-0348-8374], Filipe Andrade Bernardi¹[0000-0002-95597-5470],
Lariza Laura de Oliveira³[0000-0002-5098-172X], Afrânio Kritski²[0000-0002-5900-6007], Marcel
Koenigkam Santos³[0000-0002-7160-4691] and Domingos Alves³[0000-0002-0800-5872]

¹ Interunit Postgraduate Program in Bioengineering, University of São Paulo,
São Carlos, SP, Brazil

{marcioeloicf, filipepaulista12}@usp.br

² School of Medicine, Federal University of Rio de Janeiro, Rio de Janeiro, Brazil
{galliez77, kritskia}@gmail.com

³ Ribeirão Preto Medical School, University of São Paulo, Ribeirão Preto, São Paulo, Brazil
larizalaura@gmail.com, {marcelk46, quiron}@fmrp.usp.br

Abstract. Tuberculosis (TB), is an ancient disease that probably affects humans since pre-hominids. This disease is caused by bacteria belonging to the mycobacterium tuberculosis complex and usually affects the lungs in up to 67% of cases. In 2019, there were estimated to be over 10 million tuberculosis cases in the world, in the same year TB was between the ten leading causes of death, and the deadliest from a single infectious agent. Chest X-ray (CXR) has recently been promoted by the WHO as a tool possibly placed early in screening and triaging algorithms for TB detection. Numerous TB prevalence surveys have demonstrated that CXR is the most sensitive screening tool for pulmonary TB and that a significant proportion of people with TB are asymptomatic in the early stages of the disease. This study presents experimentation of classic convolutional neural network architectures on public CRX databases in order to create a tool applied to the diagnostic aid of TB in chest X-ray images. As result the study has an AUC ranging from 0.78 to 0.84, sensitivity from 0.76 to 0.86 and specificity from 0.58 to 0.74 depending on the network architecture. The observed performance by these metrics alone are within the range of metrics found in the literature, although there is much room for metrics improvement and bias avoiding. Also, the usage of the model in a triage use-case could be used to validate the efficiency of the model in the future.

Keywords: Tuberculosis, Chest X-Ray, Convolutional Neural Networks.

1 Introduction

1.1 Tuberculosis

Tuberculosis (TB), is an ancient disease that affects humans and probably existed in pre-hominids, and still is nowadays an important cause of death worldwide. This

disease is caused by bacteria belonging to the mycobacterium tuberculosis complex and usually affects the lungs, although other organs are affected in up to 33% of cases [1]. When properly treated, tuberculosis caused by drug-sensitive strains is curable in almost all cases. If left untreated, the disease can be fatal in 5 years in 50 to 65% of the cases. Transmission usually occurs by the aerial spread of droplets produced by patients with infectious pulmonary tuberculosis [1].

Despite the progress achieved in TB control over the past two and a half decades, with more than 50 million deaths averted globally, it is still the leading cause of death in people living with HIV, accounting for one in five deaths in the world [2]. In 2019, there were estimated to be over 10 million TB cases in the world, in the same year TB was between the ten leading causes of death, and the deadliest cause from a single infectious agent [3].

Most people who develop TB can be cured, with early diagnosis and appropriate drug treatment. Still, for many countries, the end of the disease as an epidemic and major public health problem is far from the reality. Twenty-five years ago, in 1993, WHO declared TB a global health emergency [4]. In response, the End TB Strategy has the overall goal of ending the global TB epidemic, to achieve that goal it defines the targets (2030, 2035) and milestones (2020, 2025) for the needed reductions in tuberculosis cases and deaths. The sustainable development goals include a target to end the epidemic by 2030[4].

One of these efforts supports the continued collation of the evidence and best practices for various digital health endeavors in TB prevention and care. This will make a stronger ‘investment case’ for innovative development and the essential implementation of digital health initiatives at scale [5]. Adequate triage and diagnosis are a prerequisite for the prognosis and success of any treatment and may involve several professionals and specialties. In this context, the choice of the essentially clinical mechanism that allows the measurement of the evaluator's impression becomes an important tool in helping a more accurate diagnosis [6].

1.2 Chest X-ray for detecting TB

Chest X-ray (CXR) has recently been promoted by the WHO as a tool that can be placed early in screening and triaging algorithms (see Fig. 1). A great number of prevalence surveys on TB demonstrated that CXR is the most sensitive screening tool for pulmonary TB and that TB is asymptomatic on a significant proportion of people while still in the early course of the disease [7]. When used as a triage test, CXR should be followed by further diagnostic evaluation to establish a diagnosis, it is important that any CXR abnormality consistent with TB be further evaluated with a bacteriological test [8].

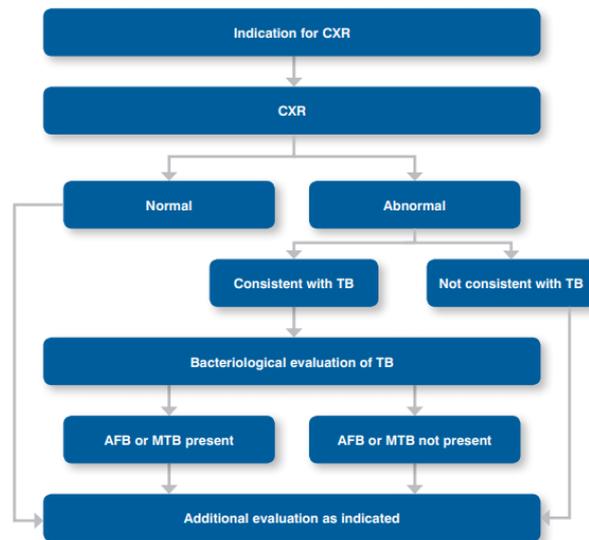


Fig. 1. Using chest radiography as a triage tool (Source: Chest radiography in tuberculosis detection – summary of current WHO recommendations and guidance on programmatic approaches, p. 11).

More than identifying active TB disease, CXR also identifies those who have inactive TB or fibrotic lesions without a history of TB treatment. Once active TB has been excluded, patients with fibrotic lesions should be followed-up, given this population is at the highest risk of developing active TB disease and/or other respiratory complications [9].

In Porto Alegre city, located in the southern region of Brazil, a projection of the impact on case detection and health system costs of alternative triage approaches for tuberculosis found that most of the triage approaches modeled without X-ray were predicted to provide no significant benefit [10]. The same study also found that adding X-ray as a triage tool for HIV-negative and unknown HIV status cases combined with appropriate triage approaches could substantially save costs over using an automated molecular test without triage, while identifying approximately the same number of cases [10].

In many high TB burden countries, it has been reported a relative lack of radiology interpretation expertise [11], this condition could result in impaired screening efficacy. There has been interest in the use of computer-aided diagnosis for the detection of pulmonary TB at chest radiography, once automated detection of disease is a cost-effective technique aiding screening evaluation [12,13]. There are already studies in the medical literature using artificial intelligence tools to evaluate CXR images, but the developed tool's availability is limited, especially in high burden countries [14].

This study presents experimentation of classic convolutional neural network architectures on public CRX databases in order to create a tool applied to the diagnostic aid of TB in chest X-ray images.

2 Methods

2.1 Google Colab

Google Colaboratory (Colab) is a project with the goal of disseminating education and research in the machine learning field [15]. Colaboratory notebooks are based on Jupyter and work as a Google Docs object, thus they can be shared as such and many users can work simultaneously on the same notebook.

Colab provides pre-configured runtimes in Python 2 or 3 with TensorFlow, Matplotlib, and Keras, essential machine learning and artificial intelligence libraries. To run the experiments, Google Colab tool provides a virtual machine with 25.51 GB RAM, GPU 1xTesla K80, having 2496 CUDA cores, CPU 1x single-core hyper-threaded Xeon Processors @2.3Ghz (No Turbo Boost), 45MB Cache.

The images were loaded into Google Drive, which can be linked to Colab for direct access to files. Google Drive provides unlimited storage disk due to a partnership between Google and the University of São Paulo (USP).

2.2 Image Datasets

PadChest. PadChest is a dataset of labeled chest x-ray images along with their associated reports. This dataset includes more than 160,000 large-scale, high-resolution images from 67,000 patients that were interpreted and reported by radiologists at Hospital San Juan (Spain) from 2009 to 2017. The images have additional data attached containing information on image acquisition and patient demography [16]. This dataset contains a total of 152 images classified with TB label.

National Institutes of Health. In 2017 the National Institutes of Health (NIH), a component of the U.S. Department of Health and Human Services, released over 100,000 anonymized chest x-ray images and their corresponding data from more than 30,000 patients, including many with advanced lung disease [17]. Other two notorious and vastly used publicly available datasets maintained by the National Institutes of Health, are from Montgomery County, Maryland, and Shenzhen, China [18] which contains respectively 58 and 336 images labeled as TB.

2.3 Network Architectures

AlexNet. In the ILSVRC classification and the localization challenge of 2012, the AlexNet architecture came on top as the winner [19]. The network architecture has 60 million parameters and 650 thousand neurons [20]. The standard settings were employed in this study: Convolutional, maxpooling, and fully connected layers, ReLU activations and the SGD optimization algorithm with a batch size of 128, momentum of 0.9, step learning annealing starting at 0.01 and reduced three times, weight decay of 0.0005, dropout layers with $p = 0.5$ design patterns, as in the original architecture article[20].

GoogLeNet. In the year of 2014, GoogLeNet was the winner of the ILSVRC detection challenge, and also came in second place on the localization challenge [21]. The standard settings were employed in this study: Convolutional, maxpooling, and fully connected layers are used, in addition, has a layer called an inception module that runs the inputs through four separate pipelines and joins them after that [22]. Also, ReLU activations, asynchronous SGD, momentum of 0.9, step learning annealing decreasing with 4% every eight epochs, a dropout layer with $p = 0.7$, as in the original architecture article [22].

ResNet. In 2015 ResNet won the classification challenge with only provided training data [23]. ResNet enables backpropagation through a shortcut between layers, this allows weights to be calculated more efficiently [24]. ResNet has a high number of layers but can be fast due to that mechanism [24]. The standard settings were employed in this study: momentum of 0.9, reduction of step learning annealing by a factor of ten every time the rate of change in error stagnates, weight decay of 0.0001 and batch normalization, as in the original architecture article [25].

2.4 Auxiliary Tools

HDF5 Dataset Generator. Functions responsible for generating the training, validation and test sets. It begins by taking a set of images and converting them to NumPy arrays, then utilizing the sklearn `train_test_split` function [26] to randomly split the images into the sets. Each set is then written to HDF5 format. HDF5 is a binary data format created by the HDF5 group [27] to store on disk numerical datasets too large to be stored in memory while facilitating easy access and computation on the rows of the datasets.

Data Augmentation Tool. Set of functions responsible for zooming, rotating and flipping the images in order to higher the generalization of the model. The Keras `ImageDataGenerator` class function [28] was utilized, with parameters set as: rotation range of 90 degrees, max zoom range of 15%, max width and height shift range of 20% and horizontal flip set as true.

2.5 Results Validation

In the experiments presented in this analysis, we choose a set of metrics based on confusion matrix [29]. Table 1 shows a confusion matrix 2x2 for a binary classifier.

Table 1. Confusion Matrix

	Actual class negative	Actual class positive
Predicted class negative	True Negative (TN)	False Negative (FN)
Predicted class positive	False Positive (FP)	True Positive (TP)

The used metrics are described in the equations below:

- Sensibility (also called true positive rate, or recall)

$$\text{Sens} = \text{TP} / (\text{TP} + \text{FN}) \quad (1)$$

- Specificity (or true negative rate)

$$\text{Spec} = \text{TN} / (\text{TN} + \text{FP}) \quad (2)$$

- Precision

$$\text{Pr} = \text{TP} / (\text{TP} + \text{FP}) \quad (3)$$

- F1-Score

$$\text{F1} = 2 * ((\text{Pr} * \text{Sens}) / (\text{Pr} + \text{Sens})) \quad (4)$$

Another metric used is the AUC - The area under the Receiver Operating Characteristic Curve (ROC) which plots the TPR (true positive rate) versus the FPR (false positive rate) [29].

3 Results

The first step was creating a new database merging the images from the Montgomery, Shenzhen, and PadChest datasets. The number of images from each dataset is described in Table 2.

Table 2. Number of images by dataset and class.

Dataset	Number of “no tb” images	Number of “yes tb” images	Total
Montgomery	80	58	138
Shinzen	326	336	662
PadChest	140	152	292
Total	546	546	1092

The image data set was split into 3 sets of images: training, validation and test sets using the HDF5 Dataset Generator. During each training the data augmentation function was responsible for preprocessing the images in order to higher the generalization of the model. The datasets sizes are described in Table 3.

Table 3. Number of images by set.

Set	Number of “no tb” images	Number of “yes tb” images
Training	446	446
Validation	50	50
Test	50	50

The first model used the AlexNet architecture by Krizhevsky et al implemented in Keras, the images were resized to $227 \times 227 \times 3$ pixels utilizing the CV2 library, in order to fit the net architecture input size. The training and validation loss/accuracy over the 75 epochs is shown in Fig. 2.

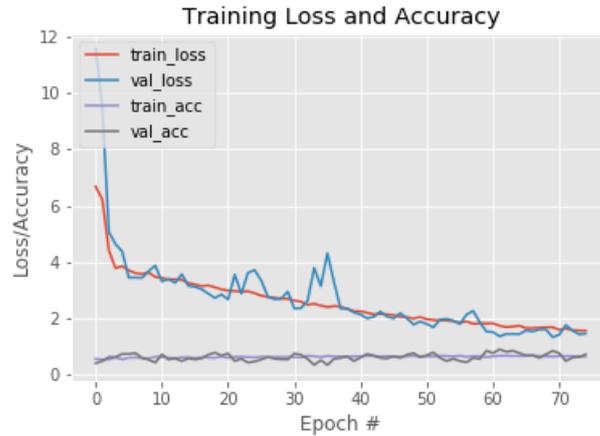


Fig. 2. AlexNet Loss and Accuracy history per Epoch.

The test confusion matrix, test performance, classification metrics and area under the ROC curve are shown in Table 4, Table 5, Table 6 and Fig. 3, respectively.

Table 4. AlexNet main classification metrics.

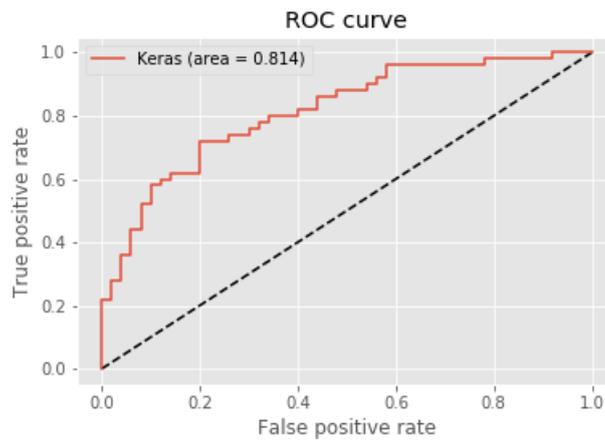
	Precision	Recall	F1-score	Support
Class: no_tb	0.68	0.86	0.76	50
Class: yes_tb	0.81	0.60	0.69	50
Accuracy	-	-	0.73	100
Macro avg	0.75	0.73	0.73	100
Weighted avg	0.75	0.73	0.73	100

Table 5. AlexNet confusion Matrix.

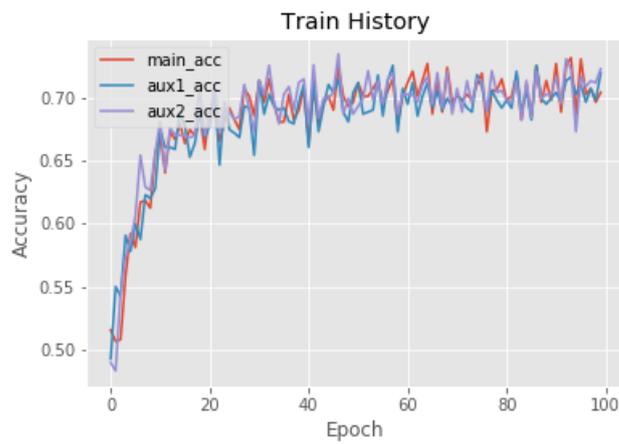
	Actual no_tb	Actual yes_tb
Predicted no_tb	43	7
Predicted yes_tb	20	30

Table 6. AlexNet model performance.

Metric	Value
Accuracy	0.73
Sensitivity	0.86
Specificity	0.60

**Fig. 3.** AlexNet ROC curve and Area Under the Curve.

Following the first experiment, the GoogLeNet architecture was implemented, the input size for this CNN model is $224 \times 224 \times 3$. The same database and datasets were utilized in the training. The training loss/accuracy over the 100 epochs is shown in Fig. 4 and Fig. 5.

**Fig. 4.** Accuracy history per Epoch.

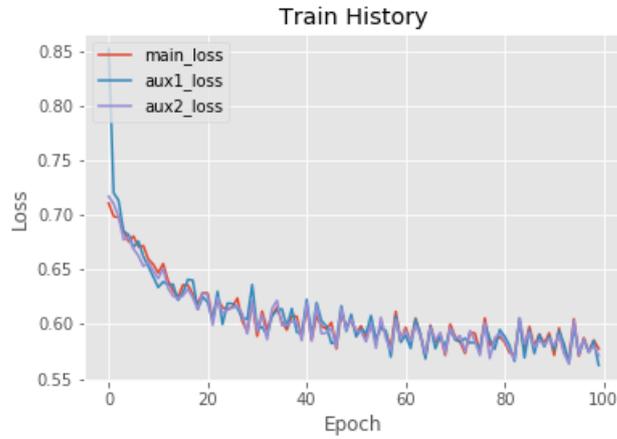


Fig. 5. GoogleNet Loss history per Epoch.

The test confusion matrix, test performance, classification metrics and area under the ROC curve are shown in Table 7, Table 8, Table 9 and Fig. 6, respectively.

Table 7. GoogleNet main classification metrics.

	Precision	Recall	F1-score	Support
Class: no_tb	0.75	0.76	0.75	50
Class: yes_tb	0.76	0.74	0.75	50
Accuracy	-	-	0.75	100
Macro avg	0.75	0.73	0.75	100
Weighted avg	0.75	0.73	0.75	100

Table 8. GoogleNet confusion Matrix.

	Actual no_tb	Actual yes_tb
Predicted no_tb	38	12
Predicted yes_tb	13	37

Table 9. GoogleNet model performance.

Metric	Value
Accuracy	0.75
Sensitivity	0.76
Specificity	0.74

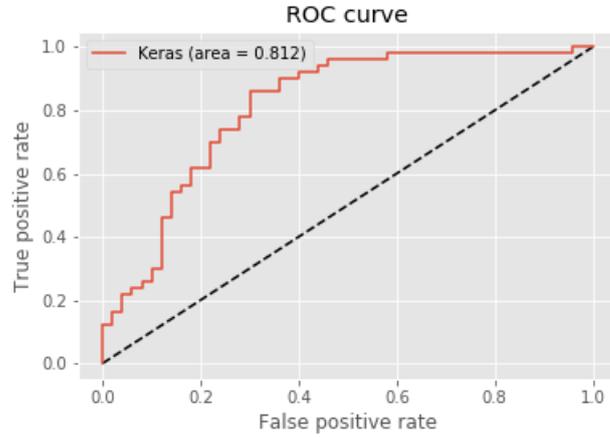


Fig. 6. GoogleNet ROC curve and Area Under the Curve.

The ResNet architecture was also implemented, trained and validated utilizing the same parameter as the other experiments. The input size for this CNN model is 224 x 224 x 3. Training and validation loss/accuracy over the 100 epochs is shown in Fig. 7.

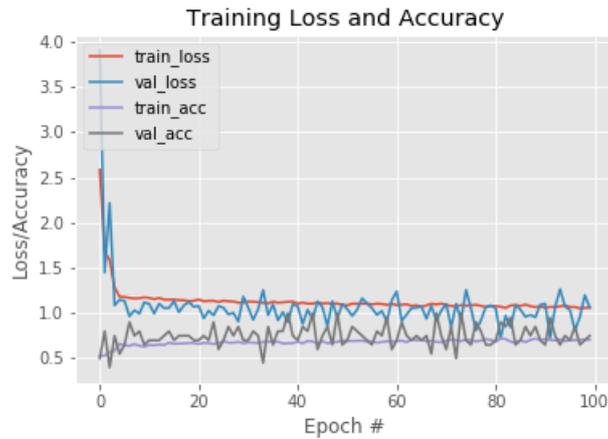


Fig. 7. ResNet Loss and Accuracy history per Epoch.

The test confusion matrix, test performance, classification metrics and area under the ROC curve are shown in Table 10, Table 11, Table 12 and Fig. 8, respectively.

Table 10. ResNet main classification metrics.

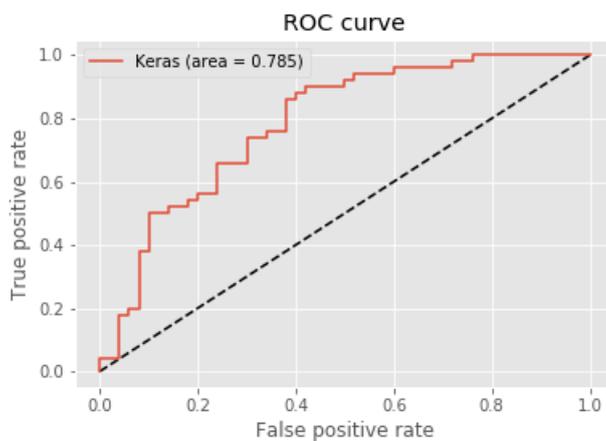
	Precision	Recall	F1-score	Support
Class: no_tb	0.64	0.76	0.70	50
Class: yes_tb	0.71	0.58	0.64	50
Accuracy	-	-	0.67	100
Macro avg	0.68	0.67	0.67	100
Weighted avg	0.68	0.67	0.67	100

Table 11. ResNet confusion Matrix.

	Actual no_tb	Actual yes_tb
Predicted no_tb	38	12
Predicted yes_tb	21	29

Table 12. ResNet model performance.

Metric	Value
Accuracy	0.67
Sensitivity	0.76
Specificity	0.58

**Fig. 8.** ResNet ROC curve and Area Under the Curve.

4 Discussion

To gather information about a higher number of articles, as well as to evaluate the quality and possible biases of each, a systematic review was chosen as a base for comparison and discussion. According to the definition presented in Harris et. al [30] this study is

classified as a Development study, that focuses on reporting methods for creating a CAD program for pulmonary TB, and includes an assessment of diagnostic accuracy.

Of the 40 studies evaluated by Harris's review, 33 reported measures of accuracy assessments, these studies had the AUC ranged from 0.78 to 0.99, sensitivity from 0.56 to 0.97, and specificity from 0.36 to 0.95. The WHO states that it is necessary for screening tests to have sensitivity greater than 0.9 and specificity greater than 0.7 [31].

This study had an AUC ranging from 0.78 to 0.84, sensitivity from 0.76 to 0.86 and specificity from 0.58 to 0.74 depending on the network architecture. The observed performance by these metrics alone are within the range of metrics found in the literature, although still far from the highest metrics obtained and did not meet the WHO standards.

The reason for such results can be speculated. One of the main possible reasons is the size of the image dataset as well as the number of datasets from different sources. A higher number of images, with high quality diagnosis, could improve the model's performance. A higher number of datasets from different sources could increase the model's generalizability.

Observing the three architecture's loss curve during training (Fig. 2, Fig 5 and Fig. 7) they seem to indicate a good fit, once training and validation loss curves both decrease to a point of stability maintaining a minimal gap between them. Although, the noisy movements on the validation line indicate an unrepresentative dataset. The performance could be improved by increasing the validation set size compared to the training set.

To further improve the model not only the metrics should be considered, but there are also many bias factors that should be avoided. The FDA (The US Food and Drug Administration) requires standards to be met for clinical use in their guidelines of CAD applied to radiology devices [32]. Those standards include a description of how CXRs were selected for training and testing, the use of images from distinct datasets for training and testing, evaluation of the model accuracy against a microbiologic reference standard and a report of the threshold score to differentiate between a positive and negative classes.

Of the requirements cited above only the first one was met by this study, which leaves open the possibility of bias.

The potential risk of bias can also be detected by applying tools for systematic reviews of diagnostic accuracy studies, the Quality Assessment of Diagnostic Accuracy Studies (QUADAS)-2 [33] is one of the approaches.

5 Conclusion

The preliminary results are between the range of metrics presented in the literature, although there is much room for improvement of metrics and bias avoiding. Also, the usage of the model in a triage use-case could be used to validate the efficiency of the model.

References

1. Fauci AS, Braunwald E, Kasper DL. Harrison. Principios de medicina interna. Vol. I . McGraw-Hill Interamericana; 2009.
2. World Health Organization (WHO). Guidelines for treatment of drug-susceptible tuberculosis and patient care, 2017 update.
3. World Health Organization. Global tuberculosis report 2019: World Health Organization; 2019.
4. World Health Organization. Global tuberculosis report 2018. World Health Organization; 2018.
5. World Health Organization. Digital health for the End TB Strategy: an agenda for action. World Health Organization; 2015.
6. ZHOU, Shang-Ming et al. Defining disease phenotypes in primary care electronic health records by a machine learning approach: a case study in identifying rheumatoid arthritis. *PLoS one*, v. 11, n. 5, p. e0154515, 2016. World Health Organization. Tuberculosis prevalence surveys: a handbook. Geneva: World Health Organization; 2011 (WHO/HTM/TB/2010.17; http://www.who.int/tb/advisory_bodies/impact_measurement_taskforce/resources_documents/thelimebook/en/, accessed 5 October 2016).
7. TB Care I. International standards for tuberculosis care, third edition. The Hague: TB CARE I; 2014 (http://www.who.int/tb/publications/ISTC_3rdEd.pdf, accessed 5 October 2016).
8. World Health Organization. Systematic screening for active tuberculosis: Principles and recommendations. Geneva: World Health Organization; 2013 (WHO/HTM/TB/2013.04; http://apps.who.int/iris/bitstream/10665/84971/1/9789241548601_eng.pdf?ua=1, accessed 27 September 2016).
9. Melendez J, Sánchez CI, Philipsen RH, et al. An automated tuberculosis screening strategy combining X-ray-based computer-aided detection and clinical information. *Sci Rep* 2016;6:25265.
10. RAHMAN, Abu AM Shazzadur et al. Modelling the impact of chest X-ray and alternative triage approaches prior to seeking a tuberculosis diagnosis. *BMC infectious diseases*, v. 19, n. 1, p. 93, 2019.
11. Antani S. Automated Detection of Lung Diseases in Chest X-Rays. A Report to the Board of Scientific Counselors. US National Library of Medicine. <https://lhncbc.nlm.nih.gov/system/files/pub9126.pdf>. Published April 2015. Accessed September 20, 2016.
12. Jaeger S, Karargyris A, Candemir S, et al. Automatic screening for tuberculosis in chest radiographs: a survey. *Quant Imaging Med Surg* 2013;3(2):89–99.
13. McAdams HP, Samei E, Dobbins J III, et al. Recent advances in chest radiography. *Radiology*. 2006;241:663–683.
14. Lakhani P, Sundaram B. Deep learning at chest radiography: automated classification of pulmonary tuberculosis by using convolutional neural networks. *Radiology*. 2017;284:574–582.
15. Colaboratory: Frequently Asked Questions, Dez. 2019, [online] Available: <https://research.google.com/colaboratory/faq.html>.
16. Bustos A, Pertusa A, Salinas JM, de la Iglesia-Vayá M. Padchest: A large chest x-ray image dataset with multi-label annotated reports. arXiv preprint arXiv:1901.07441. 2019 Jan 22.
17. Wang X, Peng Y, Lu L, Lu Z, Bagheri M, Summers RM. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition 2017* (pp. 2097-2106).

18. JAEGER, Stefan et al. Two public chest X-ray datasets for computer-aided screening of pulmonary diseases. *Quantitative imaging in medicine and surgery*, v. 4, n. 6, p. 475, 2014.
19. Imagenet large scale visual recognition challenge 2012 (ilsvrc2012). 2012. URL <http://image-net.org/challenges/LSVRC/2012/results.html>.
20. Alex Krizhevsky et al. Imagenet classification with deep convolutional neural networks. June 2017.
21. Imagenet large scale visual recognition challenge 2014 (ilsvrc2014). 2014. URL <http://image-net.org/challenges/LSVRC/2014/results>.
22. Christian Szegedy et al. Going deeper with convolutions. September 2014.
23. Imagenet large scale visual recognition challenge 2015 (ilsvrc2015). 2015. URL <http://image-net.org/challenges/LSVRC/2015/results>.
24. Øyvind Kjeldstad Grimnes. End-to-end steering angle prediction and object detection using convolutional neural networks. June 2017
25. Kaiming He et al. Deep residual learning for image recognition. December 2015.
26. Sklearn Documentation, https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html, last accessed 2020/02/03.
27. The HDF Group. Hierarchical data format version 5. <http://www.hdfgroup.org/HDF5> (cited on page 33).
28. Keras Documentation, <https://keras.io/preprocessing/image/>, last accessed 2020/02/03.
29. Maratea, A., Petrosino, A., Manzo, M.: Adjusted f-measure and kernel scaling for imbalanced data learning. *Information Sciences* 257, 331–341 (2014)
30. Harris M, Qi A, Jeagal L, Torabi N, Menzies D, Korobitsyn A, Pai M, Nathavitharana RR, Khan FA. A systematic review of the diagnostic accuracy of artificial intelligence-based computer programs to analyze chest x-rays for pulmonary tuberculosis. *PloS one*. 2019;14(9).
31. Denkinger CM, Kik SV, Cirillo DM, et al. Defining the needs for next-generation assays for tuberculosis. *J Infect Dis* 2015;211(Suppl 2):S29–38
32. Muyoyeta M, Moyo M, Kasese N, Ndhlovu M, Milimo D, Mwanza W, et al. Implementation Research to Inform the Use of Xpert MTB/RIF in Primary Health Care Facilities in High TB and HIV Settings in Resource Constrained Settings. *PLoS One*. 2015;10(6):e0126376. Epub 2015/06/02. pmid:26030301; PubMed Central PMCID: PMC4451006.
33. WHITING, Penny F. et al. QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Annals of internal medicine*, v. 155, n. 8, p. 529-536, 2011.