

The Pitfalls of Sample Selection: A Case Study on Lung Nodule Classification

Vasileios Baltatzis^{1,2}, Kyriaki-Margarita Bintsi², Loïc Le Folgoc²,
Octavio E. Martinez Manzanera¹, Sam Ellis¹, Arjun Nair³, Sujal Desai⁴,
Ben Glocker², Julia A. Schnabel^{1,5,6}

¹ School of Biomedical Engineering and Imaging Sciences, King's College London,
UK

² BioMedIA, Department of Computing, Imperial College London, UK

³ Department of Radiology, University College London, UK

⁴ The Royal Brompton & Harefield NHS Foundation Trust, London, UK

⁵ Technical University of Munich, Germany

⁶ Helmholtz Center Munich, Germany

vasileios.baltatzis@kcl.ac.uk

Abstract. Using publicly available data to determine the performance of methodological contributions is important as it facilitates reproducibility and allows scrutiny of the published results. In lung nodule classification, for example, many works report results on the publicly available LIDC dataset. In theory, this should allow a direct comparison of the performance of proposed methods and assess the impact of individual contributions. When analyzing seven recent works, however, we find that each employs a different data selection process, leading to largely varying total number of samples and ratios between benign and malignant cases. As each subset will have different characteristics with varying difficulty for classification, a direct comparison between the proposed methods is thus not always possible, nor fair. We study the particular effect of truthing when aggregating labels from multiple experts. We show that specific choices can have severe impact on the data distribution where it may be possible to achieve superior performance on one sample distribution but not on another. While we show that we can further improve on the state-of-the-art on one sample selection, we also find that on a more challenging sample selection, on the same database, the more advanced models underperform with respect to very simple baseline methods, highlighting that the selected data distribution may play an even more important role than the model architecture. This raises concerns about the validity of claimed methodological contributions. We believe the community should be aware of these pitfalls and make recommendations on how these can be avoided in future work.

1 Introduction

Lung nodule characterization is the most difficult step in the pipeline of lung cancer diagnosis according to radiologists, which can be observed by a great

inter-observer disagreement on the task [7,12]. A lung nodule is normally characterized with respect to texture, spiculation, lobulation, and its morphological appearance on a CT scan, and eventually it must be classified as either benign or malignant for patient management. The Lung imaging Reporting And Data System (Lung-RADS) [9] is a protocol that defines explicit guidelines for nodule management and follow up planning, and classifies pulmonary nodules in six categories, each of which has its own suggested follow up. Lung-RADS also integrates the PanCan Model [11], which provides a malignancy probability based on the morphology of a nodule and additional patient information. Certain diagnosis can only be made through biopsy, which, however, is invasive and not always feasible to have access to. While determining the malignancy of a nodule from its appearance on a CT scan is not a fail-proof method, it is still a very useful step of the lung cancer detection pipeline. It can have very important value to clinicians in conjunction with patient history and demographics.

Several deep learning methods have been proposed for automated nodule classification from CT. The publicly available Lung Image Database Consortium and Image Database Resource Initiative (LIDC) database [2,10] has been in the core of the majority of such efforts. The LIDC does not primarily contain pathology confirmed ground truths (besides a very small subset of cases), but rather radiologists’ annotations. Nevertheless, it is still heavily used by the research community for the task of lung nodule classification. Interestingly, there are various design choices regarding sample selection that need to be considered, which can have severe impact on the reported results.

The contributions of this paper can be summarized as follows: 1) We analyze several published works reporting results on LIDC nodule classification and examine the different assumptions such as annotation aggregations methods, removal of cases based on clinical guidelines, and data augmentation, which all can affect the resulting sample selection process; 2) Through an extensive experimental analysis, we show that the selected data distribution can affect the difficulty of the task and may play an even more important role than the model architecture; 3) We demonstrate that reproducibility and direct model comparison is virtually impossible to achieve and provide suggestions towards making this feasible in future work, while also making our data selection publicly available to promote reproducibility. We illustrate the pitfalls of sample selection with a novel methodological approach of curriculum by smoothing for lung nodule classification. Our findings and insights will be of use to the community and aid in the design of future approaches for lung nodule classification.

2 State-of-the-art in Lung Nodule Classification

The LIDC dataset contains more than 1000 scans. Each scan was reviewed by four radiologists who pinpointed lesion locations and assigned a variety of annotations including malignancy. For every nodule, each radiologist had to assign a malignancy rating from 1 (most likely benign) to 5 (most likely malignant). Nodules annotated with 3 were regarded as *indeterminate*.

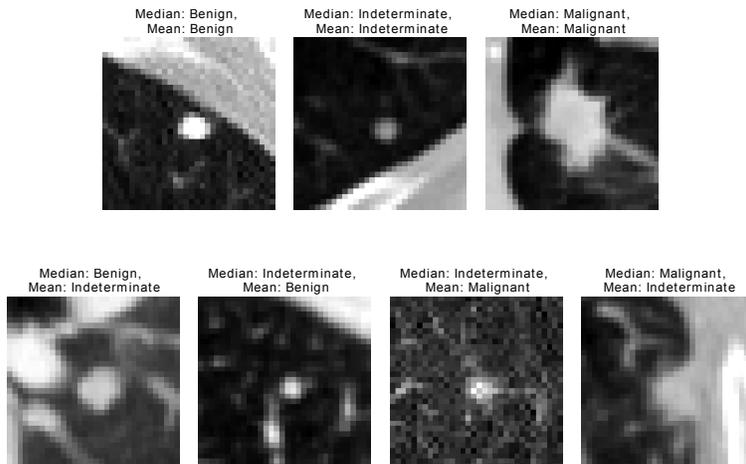


Fig. 1: Lung nodule examples from the LIDC. Top row: Nodules that have the same consensus regardless of the aggregation method used. Bottom row: Nodules that have different consensus depending on the aggregation method.

There are a number of preprocessing and data curation steps which are considered fixed when using the LIDC and almost all recent deep learning papers follow them. These include (1) retaining only nodules that have been annotated by at least three radiologists and (2) discarding nodules annotated as *indeterminate*. Subsequently, for each nodule a consensus annotation is extracted from the individual annotations through some form of aggregation or truthing (typically using mean, median, or majority voting). Example nodules from the LIDC with different consensus/aggregation combinations can be seen in Figure 1. Given these relatively straightforward steps, it may be surprising to find that every paper we studied reports largely varying numbers for benign and malignant nodules and overall cases (see Table 1). Most studies report that they follow a procedure similar to previous work, however, rarely provide the exact details about either the sample selection process or the final dataset (e.g. by publishing a list of scan series IDs). Beside the differences in absolute numbers of benign and malignant cases, the characteristics of the underlying data distribution may change significantly. One of the most important characteristics is the size of a nodule (quantified by its diameter), as it plays an essential role in malignancy classification. Another discrepancy arises from the decision to remove cases that have a slice thickness $> 2.5mm$, which is based on clinical guidelines [5]. Images with thick slices are deemed unsuitable for lung cancer screening. This step was first suggested in the LUNA16 nodule detection challenge [15] and has also been adopted by other studies [20]. One of the few works that release their pre-processed data is by Al-Shabi et al. [1].

Table 1: Overview of previous work for lung nodule classification on LIDC-IDRI in terms of nodule counts and performance. Despite all papers using the same publicly available dataset, final numbers of benign and malignant cases vary largely making a direct comparison of the methods’ performance impossible.

Method	Benign count	Malignant count	Accuracy (%)
Local-Global [1]	442	406	89.75
DeepLung [20]	554	450	90.44
Lightweight multi-CNN [14]	857	448	93.18
Interpretable hierarchical CNN [16]	3212	1040	84.20
NoduleX [3]	394	270	93.20
Multi-crop CNN [17]	880	495	87.14
Multi-task w/ margin ranking loss [8]	972	450	93.50

Here, we attempt to draw a direct comparison to their work with the dataset we have extracted from pre-processing LIDC (see Figure 2). Something like this is not feasible for the other proposed methods which do not publicly release their sample selection. In this comparison, we want to highlight the important role that the aggregation method (mean vs median) plays in determining which samples are labeled as benign and malignant. When median aggregation is used, we see that a lot more nodules have an *indeterminate* consensus (i.e. median=3) and are therefore excluded, resulting in a smaller, more balanced dataset, which is much easier to separate based on the key characteristic of nodule diameter. Specifically, median aggregation leads to 442/406 benign/malignant nodules for [1] and 376/357 benign/malignant in our replicated pipeline, respectively. In contrast, mean aggregation results in 653/484 benign/malignant for [1] and 559/451 for us. A factor leading to a discrepancy between the two samples, even when the same aggregation method is used, is the fact that cases with a slice thickness $> 2.5mm$ have been retained by [1]. These factors make reproducibility and direct comparison of methods nearly impossible.

3 Methodology

Here we present different methods and approaches, including our attempted contribution, which we considered for studying the impact of sample selection on lung nodule classification performance. We used several baselines and state-of-the-art deep learning approaches.

3.1 Diameter-based baselines

Diameter threshold The first baseline we set is not learning-based but a rather simplistic one. Specifically, given that the size of a nodule is a primary factor in determining whether a nodule is malignant or not (i.e. large nodules are most likely to be regarded by experts as malignant, while small nodules as benign) we use the provided diameter annotation in LIDC and specify a

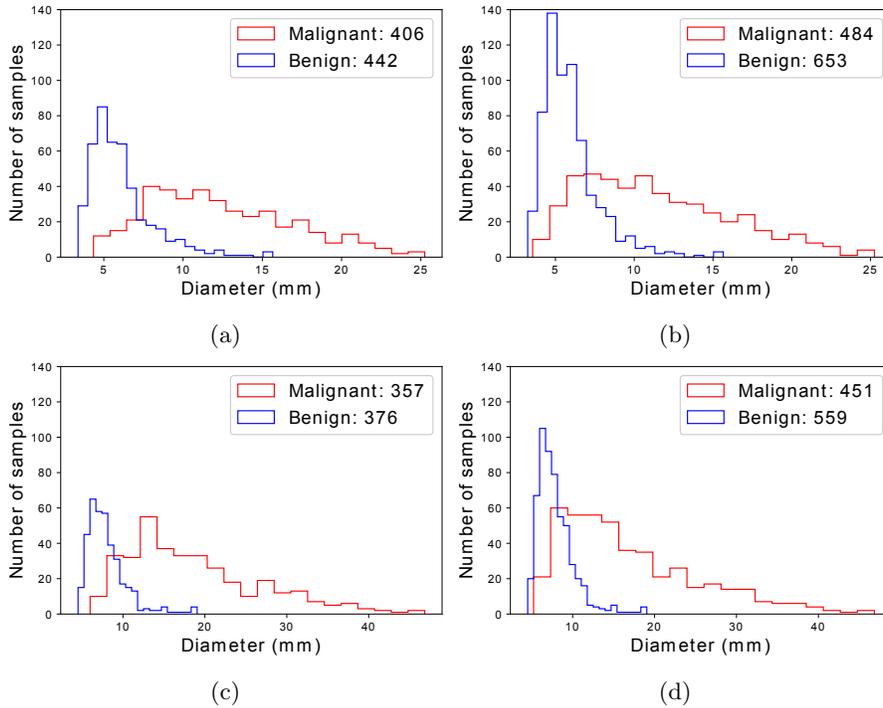


Fig. 2: Data distributions of benign and malignant samples over nodule diameter. (a) Median aggregation from [1], (b) Mean aggregation from [1], (c) Our median aggregation, (d) Our mean aggregation. Median aggregation produces fewer nodules in total (i.e. more nodules are classified as *indeterminate*) for both cases, and at the same time more balanced datasets.

threshold for classifying nodules into benign and malignant. This baseline is used as a surrogate to determine the difficulty of the classification, as the overall size difference between structures may be easily picked up by an image-based prediction model such as a convolutional neural network (CNN).

Regressed diameter threshold Another baseline that we use is similar to the previous one but with a CNN that is trained to regress the diameter through a mean squared error loss. The classification is taking place by applying the threshold determined from the first baseline on the output of the CNN instead of the annotation. Again, if this baseline works well, one may conclude that the task given a specific dataset is not very difficult.

3.2 ShallowNet

We also implement a CNN for malignancy classification (termed ShallowNet), which is a bare-bones CNN comprising of four convolutional layers with kernels

of shape 3x3 and ReLU activations, and corresponding max-pooling layers with kernels of shape 2x2, as well as a fully-connected layer with 1024 neurons at the end for the classification. This is a deliberately simplistic deep learning baseline used to compare with more complicated architectures proposed in the literature.

3.3 Local-Global

Since we have access to the sample selection of [1], it makes sense to use the state-of-the-art method on this distribution. The Local-Global network was proposed by [1] and consists of two blocks. Each block contains the following sequence: a residual sub-block [4] followed by a non-local sub-block [19] and a dropout layer. After the two blocks, there is an average pooling layer and a fully-connected layer for the classification.

3.4 Curriculum by smoothing

Finally, we propose the use of curriculum by smoothing (CBS) [18], which has shown promising results on computer vision classification tasks. CBS plays the role of our attempted methodological contribution on lung nodule classification. The main idea behind CBS is to apply a Gaussian smoothing kernel to the output of each convolutional layer of a CNN. We use $\theta \circledast x$ to denote the convolution of a kernel θ with an input x . Typically, in a CNN, a convolution operation is followed by a non-linear *activation* function as described in Equation 1:

$$z = \text{activation}(\theta_\omega \circledast x) \quad (1)$$

where θ_ω are the trainable parameters of a convolutional layer. The CBS formulation is presented in Equation 2:

$$z = \text{activation}(\theta_G \circledast (\theta_\omega \circledast x)) \quad (2)$$

where θ_G is a predefined Gaussian kernel. The Gaussian kernel is deterministic and is not trained. During the early stages of training it has an initial standard deviation σ , which is annealed as training progresses. This way, high-frequency information is suppressed in the early training steps of the CNN and is only considered at later stages of the training process.

It is important to note that while we introduce CBS here as an approach that could enhance the performance of ShallowNet or Local-Global for the task of lung nodule classification, our purpose is not to propose a novel model architecture but rather to explore whether the selected sample distribution can play a more important role than the model architecture and highlight the pitfalls that occur in such a scenario.

4 Experimental Analysis

Following from the differences in the data distributions, we move to comparing some baseline models, as well as the proposed method from [1]. In this section

we focus on two distributions to demonstrate the impact of sample selection and understand whether performance differences stem from the data or the methods. Specifically, we use the data produced with median aggregation (Figure 2a) from [1] (henceforth denoted as \mathcal{D}_1), as this is the one the authors report results for, and mean aggregation (Figure 2d) for our data (denoted as \mathcal{D}_2). We do not consider mean aggregation to be superior to median, but instead we want to study the differences in performance that are caused by this specific choice of truthing. Median aggregation leads to the two classes being more easily separated based on nodule diameter (Figures 2a,2c), even though 5-10 mm is considered the most difficult area to separate malignant from benign nodules. In both \mathcal{D}_1 and \mathcal{D}_2 , a nodule is considered benign when the consensus has a value lower than 3 and malignant when it has a value greater than 3.

CT scans with a slice thickness greater than 2.5mm are removed according to clinical guidelines [5] and every remaining scan is resampled to 1mm isotropic resolution across all three dimensions and one 32x32 mm patch is extracted along each orthogonal plane at each nodule location. The final classification result for each nodule occurs from the averaging of the individual classification of each of its three planes. Some experiments include offline data augmentation (i.e. the size of the dataset itself is increased six-fold through the addition of nodule augmentations); these augmentations are the ones suggested by [1] and include rotations, horizontal flips and Gaussian smoothing. For the proposed methodological contribution of employing CBS we choose 3x3 sized kernels, with an initial standard deviation $\sigma = 1$ of the Gaussian smoothing kernel and an annealing of 0.5 every 5 epochs based on guidelines provided by the authors of [18] and our own validation performance. All models are evaluated using 10-fold cross validation and the reported results are the average of the performance across the 10 folds. The networks are trained using the Adam optimizer [6] with learning rate 10^{-3} and binary cross-entropy loss for 50 epochs and a batch size of 256 samples. We also deploy early stopping to avoid overfitting. All experiments were conducted using PyTorch [13].

The results of the comparison can be found on Table 2. First, we show that even separating the samples based on nodule diameter (i.e. thresholding) can achieve a quite high accuracy (85.02% for \mathcal{D}_1 and 83.46% for \mathcal{D}_2). In each case, we select the threshold that maximizes training accuracy. The threshold for the two cases is quite different (7.2mm for \mathcal{D}_1 and 11.5mm for \mathcal{D}_2) because of the different aggregation methods used and also because the equivalent diameter (i.e. the diameter of the sphere having the same volume as the nodule estimated volume) is the one used in [1]. Then we use a shallow CNN (ShallowNet) to regress the nodule diameter and use a threshold (7.7mm for \mathcal{D}_1 and 11mm for \mathcal{D}_2) on that, in order to classify the nodule. If we focus on \mathcal{D}_1 , we see that a ShallowNet trained directly on malignancy can initially just outperform the diameter-based baselines (85.74%) but its performance gets better progressively when we use either CBS (86.80%) or offline augmentations (89.74%) and reaches up to 90.91% if we use both. We observe the same pattern for Local-Global [1] which starts from 89.15% when we do not use CBS or augmentations and

Table 2: Comparison of methods on the different data distribution settings. The reported results are averaged across the 10 folds. \mathcal{D}_1 is the data distribution used in [1], which has occurred from median aggregation, while \mathcal{D}_2 has been extracted from the LIDC by us using mean aggregation. We use accuracy (Acc), sensitivity (Sens) and specificity (Spec) to report the performance of each method and all the reported values are percentages (%). Even from the baselines, it is evident that \mathcal{D}_1 is an easier task to solve than \mathcal{D}_2 . All methods perform better when augmented with CBS for \mathcal{D}_1 . In \mathcal{D}_2 all configurations perform similarly to the diameter baseline, and there is no improvement from progressively increasing the complexity of the model by adding augmentations and/or CBS.

Method	\mathcal{D}_1			\mathcal{D}_2		
	Acc	Sens	Spec	Acc	Sens	Spec
Diameter threshold	85.02	90.14	80.31	83.46	69.62	94.63
CNN-regressed diameter threshold	84.43	84.23	84.61	81.58	68.95	91.77
ShallowNet	85.74	77.09	93.67	83.86	74.94	91.05
ShallowNet + CBS	86.80	78.57	94.35	82.77	71.17	92.12
ShallowNet (w/ aug)	89.74	85.96	93.21	84.35	77.38	89.98
ShallowNet (w/ aug) + CBS	90.91	89.40	92.30	82.37	73.61	89.44
Local-Global [1]	89.15	89.16	89.14	82.97	74.72	89.62
Local-Global + CBS	89.26	91.40	86.94	81.98	75.38	87.29
Local-Global (w/ aug) [1]	89.75	90.17	88.17	82.57	79.15	85.33
Local-Global (w/ aug) + CBS	90.91	90.64	91.17	81.88	70.06	91.41

eventually reaches 90.91% when we use both. The progressive gains from CBS and augmentations that are present in \mathcal{D}_1 , however, are not replicated on \mathcal{D}_2 . All the methods in that case perform very similar to the diameter-based baselines with the ShallowNet being the only one that surpasses them marginally in terms of accuracy (84.35% with augmentations).

5 Discussion

The LIDC dataset has been instrumental for the majority of recent works on lung nodule classification. Here, we take a critical look at the aspect of sample selection after discovering inconsistencies in the reported literature. We aimed to examine different factors that affect the performance of a model and thus the apparent value of its methodological contribution. Starting from the pre-processing steps that various studies have applied on the LIDC dataset, we observe that a number of different assumptions during the sample selection process can lead to very different resulting data distributions (Table 1). Such factors are the choice of the aggregation method (e.g. median or mean), in order to extract a consensus from the multiple annotations per nodule, or the removal of certain cases which are considered as unsuitable for the task due to clinical guidelines.

The aggregation method, in particular, plays a very important role. First, it is affecting the total number of nodules that are retained, since median aggregation

leads to more nodules having an *indeterminate* consensus and consequently being removed, compared to mean aggregation. It is fair to say that these nodules, which are retained in the dataset with mean aggregation, are harder examples, and therefore, the classification task that occurs from mean aggregation is more difficult. Second, the prevalence of the two classes in the dataset changes substantially, since median aggregation leads to a more balanced, and potentially more favorable for classification, dataset.

It is easy to understand that these choices change the nature of the underlying data distribution and hence, of the classification task itself. The comparison of the performance of different methods applied on different distributions is thus complex and makes the objective assessment of the value of methodological contributions difficult, which we also demonstrate experimentally. We initially devise several baselines. The first one is a simple thresholding based on the nodule diameter annotation. A size-relevant annotation is usually a core part of a lung nodule dataset, including the LIDC, and therefore this baseline can be applicable in all future studies. In the second baseline we apply a threshold on the diameter predictions that have been regressed by a neural network. This can indicate the degree of bias that a neural network has towards associating large nodules with malignancy and small ones with a benign nature. Given the very similar performance of the ShallowNet trained on malignancy prediction itself with the ShallowNet that is trained to regress the diameter, we understand that this bias is actually quite severe. It is well documented [11] that the size of the nodule is an important factor in determining whether a nodule is benign, but from a clinical perspective there are also other indications such as texture or spiculation, which do not seem to be picked up by the neural network. The aforementioned baselines can describe the difficulty of the task, and we suggest their adaptation by the research community working on lung nodule classification. Additionally, we intend to publicly release our sample selection and we urge the research community to do the same to promote reproducibility.

The core argument of our paper is epitomized when we compare the performance of all methods on the two distributions. Overall, we see that on \mathcal{D}_1 , adding data augmentation or increasing the complexity of the model (i.e. Local-Global instead of ShallowNet) consistently leads to a distinct increase in performance. The approach of using CBS during training results in a performance increase on every single method, outperforming marginally even the state-of-the-art (Local-Global w/ augmentations) on \mathcal{D}_1 . However, on \mathcal{D}_2 , all methods are bounded by the diameter threshold baseline and even CBS is not having the impact it did on \mathcal{D}_1 . This highlights the pitfalls of sample selection which may lead to incorrect conclusions about the methodological contributions. If we were to report only results on \mathcal{D}_1 , we may have concluded that CBS is beneficial for lung nodule classification, and even outperforms previous works.

6 Conclusion

In this paper we have investigated the effect of sample selection in the context of lung nodule classification using deep learning. We have investigated different factors that cause the various published studies to report completely different number of nodules, and we show experimentally that these factors explicitly affect network performance. We have demonstrated that using progressively more and more complex methods systematically improves performance on the task, if and only if the assumptions regarding the data selection process allows for it. On the other hand, if the data distribution presents a more challenging classification task, as is the case when mean aggregation for the nodule annotations is used, then model complexity or data augmentation do not offer any kind of performance boost compared to even the simplest baseline.

7 Acknowledgments

This work is funded by the King’s College London & Imperial College London EPSRC Centre for Doctoral Training in Medical Imaging (EP/L015226/1), EP-SRC grant EP/023509/1, the Wellcome/EPSCRC Centre for Medical Engineering (WT 203148/Z/16/Z), and the UKRI London Medical Imaging & Artificial Intelligence Centre for Value Based Healthcare. The Titan Xp GPU was donated by the NVIDIA Corporation.

References

1. Al-Shabi, M., Lan, B.L., Chan, W.Y., Ng, K.H., and Tan, M.: Lung nodule classification using deep Local–Global networks. *International Journal of Computer Assisted Radiology and Surgery* **14**(10), 1815–1819 (10 2019). <https://doi.org/10.1007/s11548-019-01981-7>, <https://doi.org/10.1007/s11548-019-01981-7>
2. Armato, S.G., McLennan, G., Bidaut, L., et al.: The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): A completed reference database of lung nodules on CT scans. *Medical Physics* **38**(2), 915–931 (1 2011). <https://doi.org/10.1118/1.3528204>, <http://www.ncbi.nlm.nih.gov/pubmed/21452728>
3. Causey, J.L., Zhang, J., Ma, S., et al.: Highly accurate model for prediction of lung nodule malignancy with CT scans. *Scientific Reports* **8**(1) (2018). <https://doi.org/10.1038/s41598-018-27569-w>
4. He, K., Zhang, X., Ren, S., and Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. vol. 2016-Decem, pp. 770–778 (2016). <https://doi.org/10.1109/CVPR.2016.90>, <http://image-net.org/challenges/LSVRC/2015/>
5. Kazerooni, E.A., Austin, J.H., Black, W.C., et al.: ACR-STR practice parameter for the performance and reporting of lung cancer screening thoracic computed tomography (CT): 2014 (Resolution 4). *Journal of Thoracic Imaging* **29**(5), 310–316 (2014). <https://doi.org/10.1097/RTI.000000000000097>, <https://pubmed.ncbi.nlm.nih.gov/24992501/>

6. Kingma, D.P. and Ba, J.L.: Adam: A method for stochastic optimization. In: 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings (2015)
7. Lin, H., Huang, C., Wang, W., Luo, J., Yang, X., and Liu, Y.: Measuring Interobserver Disagreement in Rating Diagnostic Characteristics of Pulmonary Nodule Using the Lung Imaging Database Consortium and Image Database Resource Initiative. *Academic Radiology* **24**(4), 401–410 (4 2017). <https://doi.org/10.1016/j.acra.2016.11.022>, <http://www.ncbi.nlm.nih.gov/pubmed/28169141>
8. Liu, L., Dou, Q., Chen, H., Qin, J., and Heng, P.A.: Multi-Task Deep Model with Margin Ranking Loss for Lung Nodule Analysis. *IEEE Transactions on Medical Imaging* **39**(3), 718–728 (3 2020). <https://doi.org/10.1109/TMI.2019.2934577>
9. McKee, B.J., Regis, S.M., McKee, A.B., Flacke, S., and Wald, C.: Performance of ACR Lung-RADS in a Clinical CT Lung Screening Program. *Journal of the American College of Radiology* **13**(2), R25–R29 (3 2016). <https://doi.org/10.1016/j.jacr.2015.12.009>, <http://www.ncbi.nlm.nih.gov/pubmed/25176499>
10. McNitt-Gray, M.F., Armato, S.G., Meyer, C.R., et al.: The Lung Image Database Consortium (LIDC) Data Collection Process for Nodule Detection and Annotation. *Academic Radiology* **14**(12), 1464–1474 (12 2007). <https://doi.org/10.1016/j.acra.2007.07.021>, <http://www.ncbi.nlm.nih.gov/pubmed/18035276>
11. McWilliams, A., Tammemagi, M.C., Mayo, J.R., et al.: Probability of cancer in pulmonary nodules detected on first screening CT. *New England Journal of Medicine* **369**(10), 910–919 (9 2013). <https://doi.org/10.1056/NEJMoa1214726>, <http://www.nejm.org/doi/10.1056/NEJMoa1214726>
12. Nair, A., Bartlett, E.C., Walsh, S.L., et al.: Variable radiological lung nodule evaluation leads to divergent management recommendations. *European Respiratory Journal* **52**(6), 1–12 (12 2018). <https://doi.org/10.1183/13993003.01359-2018>, <https://doi.org/10.1183/13993003.01359-2018>
13. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., and ...: Automatic differentiation in pytorch (2017), <https://openreview.net/forum?id=BJJsrmfCZ>
14. Sahu, P., Yu, D., Dasari, M., Hou, F., and Qin, H.: A Lightweight Multi-Section CNN for Lung Nodule Classification and Malignancy Estimation. *IEEE Journal of Biomedical and Health Informatics* **23**(3), 960–968 (5 2019). <https://doi.org/10.1109/JBHI.2018.2879834>
15. Setio, A.A.A., Traverso, A., de Bel, T., et al.: Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: The LUNA16 challenge. *Medical Image Analysis* **42**, 1–13 (12 2017). <https://doi.org/10.1016/j.media.2017.06.015>, <http://www.ncbi.nlm.nih.gov/pubmed/28732268>
16. Shen, S., Han, S.X., Aberle, D.R., Bui, A.A., and Hsu, W.: An interpretable deep hierarchical semantic convolutional neural network for lung nodule malignancy classification. *Expert Systems with Applications* **128**, 84–95 (8 2019). <https://doi.org/10.1016/j.eswa.2019.01.048>, <https://linkinghub.elsevier.com/retrieve/pii/S0957417419300545>
17. Shen, W., Zhou, M., Yang, F., et al.: Multi-crop Convolutional Neural Networks for lung nodule malignancy suspiciousness classification. *Pattern Recognition* **61**, 663–673 (1 2017). <https://doi.org/10.1016/j.patcog.2016.05.029>, <https://linkinghub.elsevier.com/retrieve/pii/S0031320316301133>

18. Sinha, S., Garg, A., and Larochelle, H.: Curriculum By Smoothing. In: Advances in Neural Information Processing Systems. vol. 33, pp. 21653–21664. Curran Associates, Inc. (2020), <https://proceedings.neurips.cc/paper/2020/file/f6a673f09493afcd8b129a0bcf1cd5bc-Paper.pdf>
19. Wang, X., Girshick, R., Gupta, A., and He, K.: Non-local Neural Networks. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 7794–7803. IEEE Computer Society (12 2018). <https://doi.org/10.1109/CVPR.2018.00813>
20. Zhu, W., Liu, C., Fan, W., and Xie, X.: DeepLung: Deep 3D Dual Path Nets for Automated Pulmonary Nodule Detection and Classification. 2018 IEEE Winter Conference on Applications of Computer Vision (WACV) (2018). <https://doi.org/10.1101/189928>, <http://arxiv.org/abs/1801.09555>