

FACULDADE DE ENGENHARIA DA UNIVERSIDADE DO PORTO

Overcoming scarce annotations through deep semi-supervised learning in cancer characterisation

Afonso José Pinheiro Oliveira Esteves Abreu



Mestrado em Engenharia Informática

Supervisor: Eduardo Rodrigues

March 10, 2025

Overcoming scarce annotations through deep semi-supervised learning in cancer characterisation

Afonso José Pinheiro Oliveira Esteves Abreu

Mestrado em Engenharia Informática

March 10, 2025

Contents

1	Introduction	1
1.1	Context	1
1.2	Problem	1
1.3	Hypothesis	2
1.4	Motivation	2
1.5	Research Questions	2
2	Literature Review	3
2.1	Eligibility criteria	3
2.2	Collection and Search strategy	4
2.3	Screening and selection process	4
2.4	Conclusion	4
	References	5

Chapter 1

Introduction

1.1 Context

Among all types of cancer, lung cancer is the one who causes the most damage in human health [4] as it is the 6th leading cause of deaths of our species, along with trachea and bronchus [12]. The reason for this remains in the fact that lung cancer is often diagnosed in a late stage of the disease. It happens that in this period, less than 10% of people survive the 5-year survival rate, which reflects the percentage of people still alive after five years of being diagnosed with this disease [23]. This imposes a major challenge in medical treatments as newer methods require earlier detection of lung cancer in order to effectively treat it.

Tissue biopsy has been the main method to identify and characterise lung cancer [5]. However, since this method requires a piece of tissue to be taken for tests, it is considered an invasive procedure, potentially leading to complications such as pneumothorax, infections, hemorrhage and damage to the tissue [26].

Computer-aided diagnosis (CAD) has gained a better reputation over the time as it can come with deep learning models able to provide a non-invasive tumour characterisation based on CT (computed tomography) scans of the human thorax [22]. This would counteract the effects of the biopsy, giving also more information about the disease more accurately.

The use of deep learning could be a major improvement in decision making support. Being able to deeply understand the characteristics of a tumour combined with the use of a non invasive procedure can help clinicians to develop targeted therapies which would then be more effective and harmless.

1.2 Problem

Training deep learning models need a lot of data for them to be reliable and robust [13]. One of the main challenges is that many existing datasets contain a significant amount of unlabelled data due to the difficulty of obtaining annotations. Especially in the case of medical images, the process is often time-consuming and costly, potentially even requiring additional examinations.

Therefore, a framework that enhances predictive performance by effectively utilizing both labelled and unlabelled data is needed.

1.3 Hypothesis

We hypothesise that a semi-supervised learning framework can effectively tackle the aforementioned challenges faced by deep learning models. As it takes into account labelled and unlabelled data, newer models could combine both types of datasets in order to have better predictive abilities [20] and help therapeutic decisions furthermore.

1.4 Motivation

The aim of better quality of life and healthcare is relevant for this paper, specially when one talks about lung cancer. Due to its complications in the early stages during diagnosis we hope we can develop a method that efficiently predicts the behavior of the disease in order to accurately indicate personalized treatments. The use of semi-supervised learning [2] can revolutionize the way we have been treating lung cancer, starting with non invasive procedures that takes into account a patient's well being.

1.5 Research Questions

This paper has the objective of creating a semi-supervised learning model in order to accurately predict the evolution of lung cancer. It makes use of supervised as well as unsupervised learning so it can potentially have more accurate results.

1. To what extent does the combination of supervised and unsupervised learning within a semi-supervised framework improve model performance for lung cancer malignancy classification, compared to using only supervised learning?
2. What is the optimal ratio of labelled to unlabelled data in the dataset(s) that will be selected for this project to achieve the best results in the semi-supervised learning model?
3. How does the number of parameters in the semi-supervised framework compare to that of a supervised learning model?
4. How does the computational time required for training and prediction vary across different combinations of models?

Chapter 2

Literature Review

Artificial Intelligence has gained significant relevance and has surely evolved rapidly through the recent years. Although it needs to be more powerful and smart to replace humans in the field of medicine, it can provide several tools to aid professionals in certain tasks which were once considered to be very time-consuming. Analyzing images and further support in medical decisions are one of the advantages we can take from Artificial Intelligence in healthcare, potentially becoming a shift of paradigm in the field.

This chapter aims to explore the recent advancements and challenges of using deep semi-supervised learning to address the shortage of annotated data in cancer characterization. Although labeled images play a crucial role in training machine learning models, these are often costly and impractical, leaving a potential spot for new deep learning models that take into account limited labeled data.

This literature review will also include several studies in the area of AI and deep semi-supervised learning regarding medical imaging. Distinct strategies and approaches used in previous investigations will be presented, each showing the obtained results, which will then be interpreted and carefully analyzed.

2.1 Eligibility criteria

For this stage, it was considered a systematic review in order to only get the pertinent studies. Eligibility criteria was established to include deep semi-supervised learning approaches in cancer characterization, which takes into account the limited annotated data. All of the researched studies were in the English language so that we could get better accessibility. Also, the papers should also have been published after the year 2020, ensuring that only the newest methods in this fast-evolving area are approached and therefore interpreted and analyzed. Besides that, it was a necessity that the papers had objective performance outcomes, so one could directly compare different models through several key aspects, such as accuracy and AUC.

To increase the number of papers obtained and also to not specify on the theme too much, other studies that focused on investigating a single technology or using deep learning models combined

with CT scans regarding other diseases were also approached, since they could turn out to be much useful for development of this thesis.

2.2 Collection and Search strategy

In order to gather relevant studies, a search strategy was conducted and it aimed to follow the PRISMA 2020 guidelines, guaranteeing diligence in the research of studies for the main problem acknowledged.

As for information sources, these included mainly PubMed, MDPI, IEEE Xplore, and Google Scholar, given their significance in areas such as healthcare and engineering. Search strategies for these databases used combined terms and keywords such as "semi-supervised learning", "cancer characterization" and "CT scans".

Another strategy that was utilized was based on the reading of a document's references. This helped since it could very much redirect to another interesting paper that could not be found in the basic search of the databases referenced above.

2.3 Screening and selection process

Regarding the data selection and screening, studies were initially picked if the title or abstract referenced the semi-supervised approach in tumor prediction for any type of cancer with the use of CT scans, always taking into account the lack of labeled data. Then, a broader set of studies that investigated certain technologies of deep learning were also picked. A clear and objective conclusion was also a key point for selection, exposing the various results obtained during the experiments. This technique ensured a somewhat quick selection process since the whole document was only read if it contained sufficient interesting information.

Although no data was extracted during this delivery, some results were carefully looked for in the studies. Such were the used algorithms, combined with their performance levels, accuracy, and AUC metrics, along with their outcomes.

2.4 Conclusion

This literature review reinforces the necessity of further research to refine approaches for enhanced interpretability and clinical applicability.

All these topics: searching, collection, screening and selecting were taken into account for the careful gathering of relevant papers and studies that could be a major benefit for this thesis, ensuring a good set of bibliographies that enriches and increases validation.

By building on the insights from these studies, future work can contribute to developing a robust solution that supports early cancer diagnosis and treatment planning.

References

- [1] Ioannis D Apostolopoulos, Nikolaos D Papathanasiou, and George S Panayiotakis. Classification of lung nodule malignancy in computed tomography imaging utilising generative adversarial networks and semi-supervised transfer learning. *Biocybernetics and Biomedical Engineering*, 41(4):1243–1257, 2021.
- [2] Yi Lin Luyang Luo Hao Chen Cheng Jin, Zhengrui Guo. Label-efficient deep learning in medical image analysis: Challenges and future directions. *Medical Image Analysis*, 2023.
- [3] Jan-Niklas Eckardt, Martin Bornhäuser, Karsten Wendt, and Jan Moritz Middeke. Semi-supervised learning in cancer diagnostics. *Frontiers in oncology*, 12:960984, 2022.
- [4] J. Ferlay, M. Ervik, F. Lam, M. Laversanne, M. Colombet, L. Mery, M. Piñeros, A. Znaor, I. Soerjomataram, and F. Bray. Global cancer observatory: Cancer today, fact sheet on all cancers, 2024. Accessed: 2024-10-09.
- [5] Z. Hou, Y. Zhan, C. Shen, W. Zhao, K. Wang, S. Yu, S. Gao, and J. Zhu. Signaling pathways driving aberrant splicing in cancer cells. *Cancer Research*, 77(6):1168–1178, 2017.
- [6] Rushi Jiao, Yichi Zhang, Le Ding, Bingsen Xue, Jicong Zhang, Rong Cai, and Cheng Jin. Learning with limited annotations: a survey on deep semi-supervised learning for medical image segmentation. *Computers in Biology and Medicine*, page 107840, 2023.
- [7] Olga Kurasova, Viktor Medvedev, Aušra Šubonienė, Gintautas Dzemyda, Aistė Gulla, Artūras Samuilis, Džiugas Jagminas, and Kęstutis Strupas. Semi-supervised learning with pseudo-labeling for pancreatic cancer detection on ct scans. In *2023 18th Iberian conference on information systems and technologies (CISTI)*, pages 1–6. IEEE, 2023.
- [8] Jinping Lao, Hongwei Lin, Haiyu Zhou, Chengchuang Lin, Zhaoliang Zheng, Gansen Zhao, and Hua Tang. Detection of high-low risk lung tumors using semi-supervised and selective labeling techniques. In *2024 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2024.
- [9] Roberto Augusto Philippi Martins and Danilo Silva. On teacher-student semi-supervised learning for chest x-ray image classification. *Anais do*, 15, 2021.
- [10] Joana Morgado, Tania Pereira, Francisco Silva, Cláudia Freitas, Eduardo Negrão, Beatriz Flor de Lima, Miguel Correia da Silva, António J Madureira, Isabel Ramos, Venceslau Hespagnol, et al. Machine learning and feature selection methods for egfr mutation status prediction in lung cancer. *Applied Sciences*, 11(7):3273, 2021.

- [11] Phuong Nguyen, Ankita Rathod, David Chapman, Smriti Prathapan, Sumeet Menon, Michael Morris, and Yelena Yesha. Active semi-supervised learning via bayesian experimental design for lung cancer classification using low dose computed tomography scans. *Applied Sciences*, 13(6):3752, 2023.
- [12] World Health Organization. The top 10 causes of death, 2024. Accessed: 2024-10-09.
- [13] C. Pinheiro, F. Silva, T. Pereira, and H.P. Oliveira. Semi-supervised approach for egfr mutation prediction on ct images. *Mathematics*, 10(4225), 2022.
- [14] Cláudia Pinheiro, Francisco Silva, Tania Pereira, and Hélder P Oliveira. Semi-supervised approach for egfr mutation prediction on ct images. *Mathematics*, 10(22):4225, 2022.
- [15] Cláudia Patrícia Ferreira Araújo Pinheiro. Ai-based cancer characterisation using semi-supervised learning algorithms. 2022.
- [16] Gil Pinheiro, Tania Pereira, Catarina Dias, Cláudia Freitas, Venceslau Hespanhol, José Luis Costa, António Cunha, and Hélder P Oliveira. Identifying relationships between imaging phenotypes and lung cancer-related mutation status: Egfr and kras. *Scientific reports*, 10(1):3625, 2020.
- [17] Feng Shi, Bojiang Chen, Qiqi Cao, Ying Wei, Qing Zhou, Rui Zhang, Yaojie Zhou, Wenjie Yang, Xiang Wang, Rongrong Fan, et al. Semi-supervised deep transfer learning for benign-malignant diagnosis of pulmonary nodules in chest ct images. *IEEE Transactions on medical imaging*, 41(4):771–781, 2021.
- [18] Francisco Silva, Tania Pereira, Joana Morgado, António Cunha, and Hélder P Oliveira. The impact of interstitial diseases patterns on lung ct segmentation. In *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 2856–2859. IEEE, 2021.
- [19] Francisco Silva, Tania Pereira, Joana Morgado, Julieta Frade, José Mendes, Cláudia Freitas, Eduardo Negrão, Beatriz Flor De Lima, Miguel Correia Da Silva, António J. Madureira, Isabel Ramos, Venceslau Hespanhol, José Luís Costa, António Cunha, and Hélder P. Oliveira. Egfr assessment in lung cancer ct images: Analysis of local and holistic regions of interest using deep unsupervised transfer learning. *IEEE Access*, 9:58667–58676, 2021.
- [20] Zahra Solatidehkordi and Imran Zualkernan. Survey on recent trends in medical image classification using semi-supervised learning. *Applied Sciences*, 12(23):12094, 2022.
- [21] V Sudharsanam, VD Vishnusriprya, Yagnavajjula Likhitha, R Thenila, et al. Detection of covid-19 on lung ct images using semi supervised learning. In *2021 2nd International Conference on Smart Electronics and Communication (ICOSEC)*, pages 973–977. IEEE, 2021.
- [22] M.A. Thanoon, M.A. Zulkifley, M.A.A. Mohd Zainuri, and S.R. Abdani. A review of deep learning techniques for lung cancer screening and diagnosis based on ct images. *Diagnostics*, 13(16):2617, 2023.
- [23] T. Baird A.M. et al. The Health Policy Partnership, Albrecht. Lung cancer in europe: the way forward, 2022. Accessed: 2024-10-09.

- [24] Guotai Wang, Shuwei Zhai, Giovanni Lasio, Baoshe Zhang, Byong Yi, Shifeng Chen, Thomas J Macvittie, Dimitris Metaxas, Jinghao Zhou, and Shaoting Zhang. Semi-supervised segmentation of radiation-induced pulmonary fibrosis from lung ct scans with multi-scale guided dense attention. *IEEE transactions on medical imaging*, 41(3):531–542, 2021.
- [25] Lulu Wang. Deep learning techniques to diagnose lung cancer. *Cancers*, 14(22):5569, 2022.
- [26] M.R. Wilkins and K.A. Paschke. Clinical practice. cystic fibrosis respiratory infections: optimizing treatment in the era of cftr modulators. *PubMed*, 364(9428):1195–1206, 2011.
- [27] Yutong Xie, Jianpeng Zhang, and Yong Xia. Semi-supervised adversarial model for benign–malignant lung nodule classification on chest ct. *Medical image analysis*, 57:237–248, 2019.
- [28] Xiangli Yang, Zixing Song, Irwin King, and Zenglin Xu. A survey on deep semi-supervised learning. *IEEE Transactions on Knowledge and Data Engineering*, 35(9):8934–8954, 2022.
- [29] Guojin Zhang, Yuntai Cao, Jing Zhang, Jialiang Ren, Zhiyong Zhao, Xiaodi Zhang, Shenglin Li, Liangna Deng, and Junlin Zhou. Predicting egfr mutation status in lung adenocarcinoma: development and validation of a computed tomography-based radiomics signature. *American journal of cancer research*, 11(2):546, 2021.