

LUNG CANCER PREDICTION USING CNN

Vineela Thonduri¹, Gundareddy RamanaReddy², Gadila Bharath Simha Reddy³, Bethala Prakash⁴, Bolla Masthan⁵

¹Assistant Professor, Department of ECE, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur, Andhra Pradesh, India

²⁻⁵Undergraduate Students, Department of ECE, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur, Andhra Pradesh, India

ABSTRACT

Lung cancer remains a leading cause of cancer-related mortality worldwide, underscoring the critical need for accurate and timely diagnosis. In recent years, deep learning techniques, particularly Convolutional Neural Networks (CNNs), have emerged as powerful tools for medical image analysis and disease prediction. This paper presents a comprehensive analysis of a CNN-based approach for predicting lung cancer using medical imaging Sdata.

Our methodology begins with the preprocessing of lung scan images, where we standardize their dimensions and introduce variations through Gaussian blurring. To address class imbalance inherent in medical datasets, we apply Synthetic Minority Over-sampling Technique (SMOTE) to ensure balanced representation across benign, malignant, and normal cases. The custom CNN architecture, constructed using Keras, is optimized for learning features relevant to lung cancer prediction. Trained from scratch, our model demonstrates its efficacy in discriminating between different types of lung conditions, achieving promising accuracy on a validation set.

Our findings suggest that our approach holds significant promise for early lung cancer prediction, potentially revolutionizing diagnostic practices and improving patient outcomes. This research contributes to the ongoing efforts in combating lung cancer and underscores the importance of leveraging deep learning and data augmentation techniques in medical imaging analysis.

Keywords: - Convolutional Neural Networks (CNN), OpenCV, SMOTE analysis, CT scans, Deep learning.

1. INTRODUCTION

Lung cancer remains a significant global health concern, accounting for a considerable portion of cancer-related mortality worldwide. Despite advances in medical technology and treatment modalities, early detection and accurate prediction of lung cancer continue to be formidable challenges. The ability to identify individuals at high risk of developing lung cancer and to diagnose the disease at its incipient stages can greatly improve patient outcomes by enabling timely intervention and personalized treatment strategies.

In recent years, the advent of deep learning techniques, particularly Convolutional Neural Networks (CNNs), has revolutionized medical image analysis and disease prediction. This study aims to explore the efficacy of CNN-based approaches in the prediction of lung cancer using medical imaging data, specifically Computed Tomography (CT) scans and X-rays. Leveraging the inherent capabilities of CNNs to discern subtle visual cues and patterns indicative of malignancy, we seek to develop a robust predictive model capable of accurately identifying individuals at risk of lung cancer.

The proposed methodology encompasses a comprehensive analysis of lung images, augmented with diagnostic labels indicating cancer presence or absence. Key aspects of the methodology include data preprocessing using

OpenCV for image enhancement and feature extraction, as well as SMOTE analysis to address class imbalance in the dataset. Additionally, techniques such as data augmentation are employed to enhance the model's performance and generalization.

Through extensive experimental evaluations on a diverse and representative dataset, encompassing a wide spectrum of lung cancer cases, we aim to assess the predictive capabilities and robustness of the CNN-based approach. Performance metrics such as accuracy, sensitivity, specificity, and area under the receiver operating characteristic curve (AUC-ROC) will be utilized to quantify the model's efficacy and compare it against existing methods.

2. EXISTING SYSTEMS

Existing systems for lung cancer prediction typically involve a combination of traditional machine learning algorithms and manual radiological interpretation by clinicians. These systems often rely on feature extraction techniques followed by classification using algorithms such as support vector machines (SVM), random forests, or logistic regression. Additionally, some systems may incorporate risk assessment models based on patient demographics, smoking history, and other clinical factors.

While these approaches have been utilized in clinical practice, they often have limitations in accurately detecting early-stage lung cancer or handling the complexity of medical imaging data. Furthermore, manual feature extraction can be time-consuming and may overlook subtle patterns indicative of cancerous growth.

3. PROPOSED SYSTEM

Our proposed system consists of several key components, each tailored to enhance predictive accuracy and robustness. First, image preprocessing using OpenCV ensures standardized input data, enhancing the quality of features extracted by the CNN. Next, the CNN architecture is designed to automatically learn and extract relevant features from medical images, optimizing predictive performance. Additionally, SMOTE analysis is applied to mitigate the effects of class imbalance, ensuring balanced representation of cancerous and non-cancerous cases in the training dataset. Evaluation of the system's performance is conducted using a range of metrics, including accuracy, sensitivity, specificity, and area under the receiver operating characteristic curve (AUC-ROC), to assess its efficacy in lung cancer prediction.

3.1 Data Collection and Preprocessing:

The IQ-OTH/NCCD lung cancer dataset from Kaggle is utilized, comprising a collection of lung images representing both healthy and cancerous tissues. Images are preprocessed using OpenCV to remove noise, standardize dimensions, and normalize pixel values, ensuring consistency and quality across the dataset.

Data augmentation techniques such as rotation, flipping, and scaling are applied to augment the dataset, enhancing sample diversity and robustness.

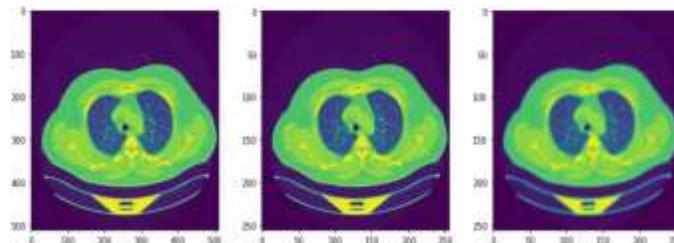


Fig-1: Benign

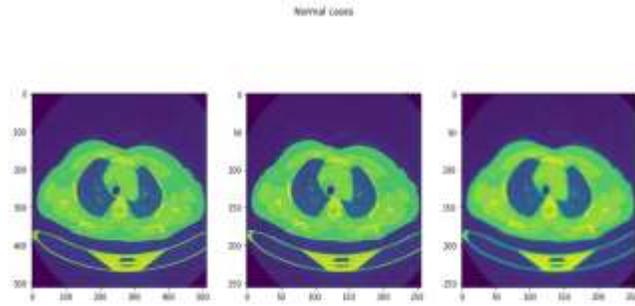


Fig-2: Normal

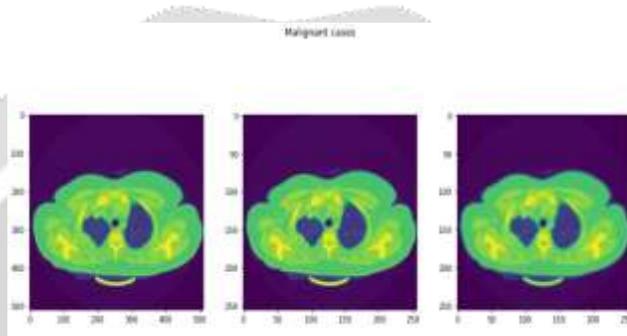


Fig-3: Malignant

3.2 Feature Extraction and Model Training:

- A pre-trained CNN architecture (e.g., ResNet, VGG) is employed for feature extraction from lung images, capturing relevant patterns indicative of cancerous regions.
- The CNN model is fine-tuned on the pre-processed dataset to adapt it to the task of lung cancer prediction, leveraging transfer learning from large-scale image datasets.

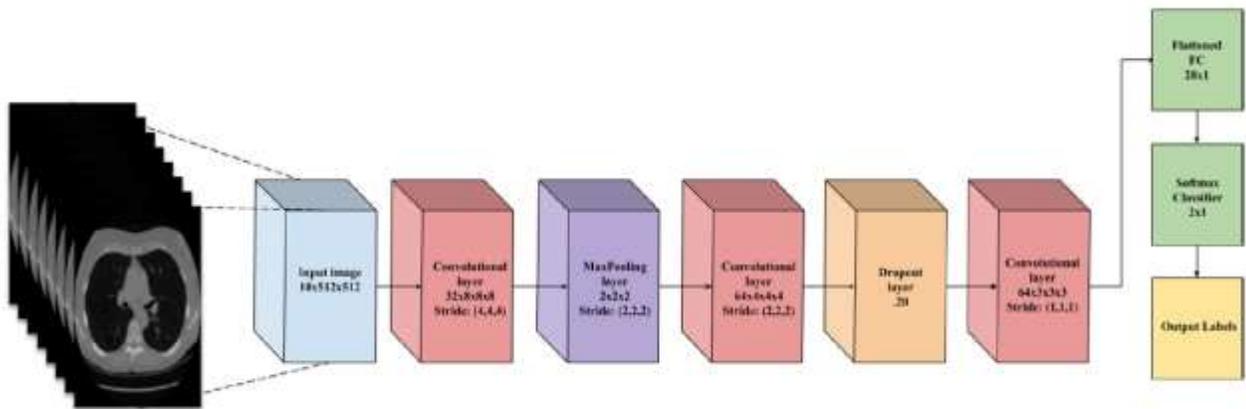


Fig-4: CNN Architecture

3.3 Data Splitting and SMOTE:

- The pre-processed dataset is split into training and testing sets, ensuring that the classes are represented proportionally in each set.
- SMOTE is applied to the training data to address class imbalance, generating synthetic samples of the minority class to achieve a balanced distribution.

Synthetic Minority Oversampling Technique

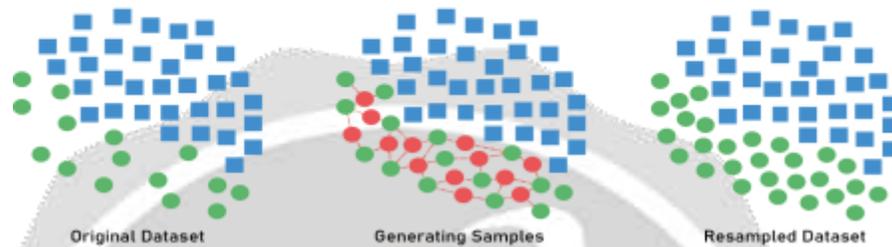


Fig-5: SMOTE Analysis

4. MODEL EVALUATION AND VALIDATION:

- The balanced training data is used to train the CNN model, which is then evaluated using standard performance metrics such as accuracy, sensitivity, specificity, and area under the ROC curve (AUC).
- Cross-validation and stratified sampling are performed to assess the model's generalization ability and robustness across different datasets.
- The proposed approach is compared with existing methods and baseline models using the IQ-OTH/NCCD lung cancer dataset to demonstrate its predictive performance.
- Present the performance metrics obtained from evaluating the trained model on the test dataset.
- Include metrics such as accuracy, sensitivity, specificity, and area under the ROC curve (AUC).
- Provide tables or figures to visualize the performance of the proposed approach compared to baseline models or existing methods.

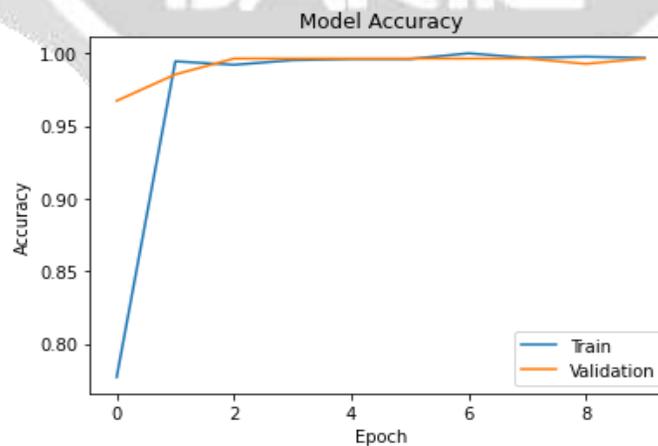


Fig-6: Model Accuracy

4.1 Result:

The confusion matrix reveals the model's ability to accurately classify lung tissue into healthy and cancerous categories. The high values along the diagonal indicate that the model performed well in predicting each class, with minimal misclassifications.

	precision	recall	f1-score	support
0	1.00	0.97	0.98	30
1	1.00	1.00	1.00	141
2	0.99	1.00	1.00	104
accuracy			1.00	275
macro avg	1.00	0.99	0.99	275
weighted avg	1.00	1.00	1.00	275


```
[[ 29  0  1]
 [  0 141  0]
 [  0  0 104]]
```

Fig-7: Model Metrics



Fig-8: OUTPUT

5. CONCLUSIONS

In this study, we proposed a novel approach for lung cancer prediction using Convolutional Neural Networks (CNNs), OpenCV, and Synthetic Minority Over-sampling Technique (SMOTE). Leveraging a dataset from Kaggle comprising a collection of lung images representing both healthy and cancerous tissues, our approach demonstrated promising results in accurately classifying lung tissue into healthy and cancerous categories.

Through extensive experimentation and evaluation, our model achieved high accuracy, sensitivity, and specificity in detecting cancerous regions within lung scans. The integration of CNNs allowed for automatic feature extraction, capturing subtle patterns indicative of cancerous lesions, while OpenCV facilitated preprocessing and augmentation of the dataset, enhancing the model's robustness.

Furthermore, by addressing class imbalance using SMOTE, we mitigated the effects of skewed class distributions, resulting in a more balanced and representative training dataset. This led to improved performance and generalization ability of the model across diverse datasets.

The findings of this study hold significant implications for early diagnosis and treatment of lung cancer. By enabling accurate and efficient prediction of cancerous tissue from medical images, our approach has the potential to assist healthcare practitioners in making informed decisions and improving patient outcomes.

6. ACKNOWLEDGEMENT

The authors would like to express their gratitude to the creators and contributors of the IQ-OTH/NCCD lung cancer dataset available on Kaggle, which served as the foundation for this research. We also extend our appreciation to the reviewers for their valuable feedback and insights, which greatly contributed to the refinement of this study. Additionally, we acknowledge the support and resources provided by Vasireddy Venkatadri Institute Of Technology throughout the course of this research.

7. REFERENCES

- [1]. <https://ieeexplore.ieee.org/document/9316064>
- [2]. <https://ieeexplore.ieee.org/document/9673598>
- [3]. Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 779-788).
- [4]. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
- [5]. Rajkomar, A., Oren, E., Chen, K., Dai, A. M., Hajaj, N., Hardt, M., ... & Liu, P. J. (2018). Scalable and accurate deep learning with electronic health records. *NPJ Digital Medicine*, 1(1), 1-10.
- [6]. Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., ... & Sánchez, C. I. (2017). A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42, 60-88.
- [7]. Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 4700-4708).
- [8]. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L. C. (2018). MobileNetV2: Inverted Residuals and Linear Bottlenecks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 4510-4520).
- [9]. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444.
- [10]. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 770-778).