

Machine Learning and Genetic Algorithm for Lung Cancer Imaging Analysis

Abstract: The study explores the use of Convolutional Neural Networks (CNNs) and Random Forest models for lung cancer detection. Using two datasets, the IQ-OTHNCCD and LC25000 Lung Datasets, CNN achieved 96.67% accuracy and 100% precision, respectively. Random Forest showed moderate performance, but GA optimization and K-fold cross validation improved its accuracy. Future research will explore further optimization strategies and clinical data integration for improved diagnostic capabilities.

- ⇒ গবেষণাটি ফুসফুসের ক্যান্সার সনাক্তকরণের জন্য কনভোলিউশনাল নিউরাল নেটওয়ার্ক (সিএনএন) এবং র্যান্ডম ফরেস্ট মডেলের ব্যবহার অন্বেষণ করে। দুটি ডেটাসেট ব্যবহার করে, IQ-OTHNCCD এবং LC25000 ফুসফুসের ডেটাসেট, CNN যথাক্রমে 96.67% নির্ভুলতা এবং 100% নির্ভুলতা অর্জন করেছে। র্যান্ডম ফরেস্ট মাঝারি পারফরম্যান্স দেখিয়েছে, কিন্তু GA অপটিমাইজেশন এবং কে-ফোল্ড ক্রস বৈধতা এর নির্ভুলতা উন্নত করেছে। ভবিষ্যত গবেষণা উন্নত ডায়গনস্টিক ক্ষমতার জন্য আরও অপটিমাইজেশন কৌশল এবং ক্লিনিকাল ডেটা ইন্টিগ্রেশন অন্বেষণ করবে।

Introduction:

Lung cancer is a global health challenge, affecting an estimated 2.2 million lives in 2020 alone. Early detection is crucial for improving patient prognosis and survival rates. Traditional methods, such as chest X-rays and CT scans, have limitations due to human error and inter-reader variability. Recent advancements in imaging technologies and AI have increased lung cancer diagnosis efficiency and accuracy. Deep learning and machine learning techniques have emerged as powerful tools for improving lung cancer imaging analysis. AI methods, such as ML and DL, can address the complexities of imaging techniques like PET, MRI, CT, and ultrasound. Random Forest, an ensemble learning method, has shown promise in medical image classification. Convolutional Neural Networks (CNNs) and Transformers are at the forefront of AI-driven breakthroughs in medical imaging. These approaches push the boundaries of early detection of lung cancer, providing clinicians with robust tools for more accurate diagnoses.

- ⇒ ফুসফুসের ক্যান্সার একটি বিশ্বব্যাপী স্বাস্থ্য চ্যালেঞ্জ, যা শুধুমাত্র 2020 সালে আনুমানিক 2.2 মিলিয়ন জীবনকে প্রভাবিত করে। রোগীর পূর্বাভাস এবং বেঁচে থাকার হার উন্নত করার জন্য প্রাথমিক সনাক্তকরণ অত্যন্ত গুরুত্বপূর্ণ। প্রথাগত পদ্ধতি, যেমন বুকের এক্স-রে এবং সিটি স্ক্যান, মানুষের ত্রুটি এবং আন্তঃপাঠক পরিবর্তনশীলতার কারণে সীমাবদ্ধতা রয়েছে। ইমেজিং প্রযুক্তি এবং এআই-এর সাম্প্রতিক অগ্রগতি ফুসফুসের ক্যান্সার নির্ণয়ের দক্ষতা এবং নির্ভুলতা বাড়িয়েছে। ডিপ লার্নিং এবং মেশিন লার্নিং কৌশলগুলি ফুসফুসের ক্যান্সার ইমেজিং বিশ্লেষণের উন্নতির জন্য শক্তিশালী সরঞ্জাম হিসাবে আবির্ভূত হয়েছে। এআই পদ্ধতি, যেমন ML এবং DL, PET, MRI, CT এবং আল্ট্রাসাউন্ডের মতো ইমেজিং কৌশলগুলির জটিলতাগুলিকে মোকাবেলা করতে পারে। র্যান্ডম ফরেস্ট, একটি সমন্বিত শিক্ষা পদ্ধতি, চিকিৎসা চিত্র শ্রেণীবিভাগে প্রতিশ্রুতি দেখিয়েছে। কনভোলিউশনাল নিউরাল নেটওয়ার্ক (সিএনএন) এবং ট্রান্সফরমারগুলি মেডিকেল ইমেজিংয়ে এআই-চালিত সাফল্যের অগ্রভাগে রয়েছে। এই পদ্ধতিগুলি ফুসফুসের ক্যান্সারের প্রাথমিক সনাক্তকরণের সীমানাকে ঠেলে দেয়, আরও সঠিক নির্ণয়ের জন্য চিকিত্সকদের শক্তিশালী সরঞ্জাম সরবরাহ করে।

Literature review: Lung cancer diagnosis is crucial for survival rates, and current methods have limitations. Artificial intelligence and deep learning have introduced new possibilities, particularly through Convolutional Neural Networks (CNNs). CNNs can outperform human radiologists in low-dose CT scans, reducing false positives and missed diagnoses. However, limitations include overfitting, limited comparative analysis, and misclassification. Recent advancements in deep learning extend CNN applications beyond lung nodule detection.

⇒ ফুসফুসের ক্যান্সার নির্ণয় বেঁচে থাকার হারের জন্য অত্যন্ত গুরুত্বপূর্ণ এবং বর্তমান পদ্ধতির সীমাবদ্ধতা রয়েছে। কৃত্রিম বুদ্ধিমত্তা এবং গভীর শিক্ষা নতুন সম্ভাবনার সূচনা করেছে, বিশেষ করে কনভোল্যুশনাল নিউরাল নেটওয়ার্ক (সিএনএন) এর মাধ্যমে। সিএনএন কম ডোজ সিটি স্ক্যানে মানব রেডিওলজিস্টদের ছাড়িয়ে যেতে পারে, মিথ্যা ইতিবাচক এবং মিস ডায়াগনোসিস কমাতে পারে। যাইহোক, সীমাবদ্ধতার মধ্যে ওভারফিটিং, সীমিত তুলনামূলক বিশ্লেষণ এবং ভুল শ্রেণীকরণ অন্তর্ভুক্ত। গভীর শিক্ষার সাম্প্রতিক অগ্রগতি ফুসফুসের নডিউল সনাক্তকরণের বাইরে সিএনএন অ্যাপ্লিকেশনগুলিকে প্রসারিত করেছে।

Methodology:

Inclusion Criteria	Exclusion Criteria
Articles published in peer-reviewed journals or presented at reputable conferences. Articles published in peer-reviewed journals or presented at reputable conferences. - Articles written in English or with English translations available.	Articles published in peer-reviewed journals or presented at reputable conferences. Studies unrelated to deep learning or lung cancer imaging analysis. - Articles written in languages other than English without available translations.

Traditional Methods for Lung Cancer Diagnosis:

Lung cancer diagnosis traditionally relies on chest X-rays and low-dose CT scans, which have limited sensitivity for early-stage lung cancer. Low-dose CT scans have improved sensitivity but are time-consuming and laborious, leading to missed diagnoses. This highlights the need for more objective, efficient tools.

⇒ ফুসফুসের ক্যান্সার নির্ণয় ঐতিহ্যগতভাবে বুকের এক্স-রে এবং কম ডোজ সিটি স্ক্যানের উপর নির্ভর করে, যার প্রাথমিক পর্যায়ে ফুসফুসের ক্যান্সারের জন্য সীমিত সংবেদনশীলতা রয়েছে। কম-ডোজ সিটি স্ক্যানগুলি সংবেদনশীলতা উন্নত করেছে কিন্তু সময়সাপেক্ষ এবং শ্রমসাধ্য, যার ফলে মিস ডায়াগনোসিস হয়। এটি আরও উদ্দেশ্যমূলক, দক্ষ সরঞ্জামগুলির প্রয়োজনীয়তা তুলে ধরে।

Algorithms:

Reinforcement Learning (RL) involves agents interacting with environments to learn optimal policies for sequential decision making. RL agents are typically categorized as model-based or model-less, with model-based agents relying on models of the environment to predict subsequent states and rewards.

- ⇒ রিইনফোর্সমেন্ট লার্নিং (RL) ক্রমিক সিদ্ধান্ত গ্রহণের জন্য সর্বোত্তম নীতিগুলি শিখতে পরিবেশের সাথে যোগাযোগকারী এজেন্টদের জড়িত করে। আরএল এজেন্টদের সাধারণত মডেল-ভিত্তিক বা মডেল-লেস হিসাবে শ্রেণীবদ্ধ করা হয়, মডেল-ভিত্তিক এজেন্টরা পরবর্তী অবস্থা এবং পুরস্কারের পূর্বাভাস দিতে পরিবেশের মডেলের উপর নির্ভর করে।

Datasets:

This study analyzes the IQ-OTH/NCCD lung cancer dataset, used for deep learning to identify and classify lung cancers. It also examines the LC25000 lung cancer dataset, a collection of preprocessed images, for in-depth examination and enhancement of models for accurate disease detection. Both datasets provide diverse images for comprehensive assessment and improvement.

- ⇒ এই অধ্যয়নটি IQ-OTH/NCCD ফুসফুসের ক্যান্সার ডেটাসেট বিশ্লেষণ করে, যা ফুসফুসের ক্যান্সার সনাক্ত এবং শ্রেণীবদ্ধ করতে গভীর শিক্ষার জন্য ব্যবহৃত হয়। এটি LC25000 ফুসফুসের ক্যান্সার ডেটাসেট, প্রি-প্রসেসড ইমেজগুলির একটি সংগ্রহ, গভীরভাবে পরীক্ষা এবং সঠিক রোগ সনাক্তকরণের জন্য মডেলগুলির উন্নতির জন্য পরীক্ষা করে। উভয় ডেটাসেটই ব্যাপক মূল্যায়ন এবং উন্নতির জন্য বিভিন্ন চিত্র প্রদান করে।

Results Analysis:

The study analyzed the performance of Convolutional Neural Networks (CNNs) and Random Forest models on the IQ-OTHNCCD lung cancer dataset and the LC25000 lung cancer dataset. The Random Forest model outperformed CNNs in terms of accuracy, achieving more accurate results. Data augmentation was applied to increase variability and improve model generalization. The Random Forest Classifier was employed as the classification model, aggregating multiple decision trees for robust predictions. Grid Search CV was used to optimize the model hyper parameters. The results showed that the Random Forest model outperformed CNNs in terms of accuracy, F1-score, recall, and precision. Future work could explore hybrid approaches that combine the strengths of both CNNs and Random Forests or investigate other advanced models to improve classification accuracy and generalizability.

- ⇒ গবেষণাটি IQ-OTHNCCD ফুসফুসের ক্যান্সার ডেটাসেট এবং LC25000 ফুসফুসের ক্যান্সার ডেটাসেটে কনভোলিউশনাল নিউরাল নেটওয়ার্ক (CNNs) এবং র্যান্ডম ফরেস্ট মডেলগুলির কর্মক্ষমতা বিশ্লেষণ করেছে। র্যান্ডম ফরেস্ট মডেল নির্ভুলতার পরিপ্রেক্ষিতে সিএনএন-কে ছাড়িয়ে গেছে, আরও সঠিক ফলাফল অর্জন করেছে। ডেটা পরিবর্তনশীলতা বৃদ্ধি এবং মডেল সাধারণীকরণ উন্নত করতে বর্ধন প্রয়োগ করা হয়েছিল। র্যান্ডম ফরেস্ট ক্লাসিফায়ারকে শ্রেণীবিন্যাস মডেল হিসাবে নিযুক্ত করা হয়েছিল, শক্তিশালী করার জন্য একাধিক সিদ্ধান্ত গাছকে একত্রিত করে ভবিষ্যদ্বাণী মডেল হাইপার প্যারামিটার অপটিমাইজ করতে গ্রিড সার্চ সিভি ব্যবহার করা হয়েছিল। ফলাফলগুলি দেখায় যে র্যান্ডম ফরেস্ট মডেলটি নির্ভুলতা, এফ1-স্কোর, রিকল এবং নির্ভুলতার ক্ষেত্রে সিএনএনকে ছাড়িয়ে গেছে। ভবিষ্যত কাজ হাইব্রিড পদ্ধতির অন্বেষণ করতে পারে যা CNN এবং র্যান্ডম ফরেস্ট উভয়ের শক্তিকে একত্রিত করে বা শ্রেণীবিভাগের নির্ভুলতা এবং সাধারণীকরণের উন্নতির জন্য অন্যান্য উন্নত মডেলগুলি তদন্ত করে।

Discussion:

The study compares the performance of CNNs and Random Forest models in lung cancer imaging analysis on two datasets: the IQ-OTHNCCD and LC25000 Lung Datasets. Random Forest, when enhanced with optimization techniques like Genetic Algorithms (GA) and K-fold cross-validation, offers exceptional accuracy in classifying lung cancer subtypes. Its adaptability to complex datasets outperforms CNN in experiments, particularly on the LC25000 dataset. However, CNN's performance on the LC25000 dataset was lower, suggesting potential limitations in deep learning models for certain medical data types. Gene Transformer, a deep learning framework that integrates CNN and attention mechanisms, is particularly effective at handling high-dimensional data, such as gene expression or medical imaging datasets. It also addresses the interpretability gap in Random Forest models, allowing researchers to gain insights into specific genes or imaging features driving classification decisions. Future studies could integrate clinical and histopathological data with deep learning models for more comprehensive cancer analysis.

- ⇒ গবেষণাটি দুটি ডেটাসেটে ফুসফুসের ক্যান্সার ইমেজিং বিশ্লেষণে সিএনএন এবং র্যান্ডম ফরেস্ট মডেলের কর্মক্ষমতা তুলনা করে: IQ-OTHNCCD এবং LC25000 ফুসফুসের ডেটাসেট। র্যান্ডম ফরেস্ট, যখন জেনেটিক অ্যালগরিদম (GA) এবং কে-ফোল্ড ক্রস-ভ্যালিডেশনের মতো অপটিমাইজেশন কৌশলগুলির সাথে উন্নত করা হয়, তখন ফুসফুসের ক্যান্সারের সাবটাইপগুলিকে শ্রেণীবদ্ধ করার ক্ষেত্রে ব্যতিক্রমী নির্ভুলতা প্রদান করে। জটিল ডেটাসেটের সাথে এর অভিযোজনযোগ্যতা পরীক্ষায় CNN কে ছাড়িয়ে যায়, বিশেষ করে LC25000 ডেটাসেটে। যাইহোক, LC25000 ডেটাসেটে CNN-এর কর্মক্ষমতা কম ছিল, যা কিছু মেডিকেল ডেটা প্রকারের জন্য গভীর শিক্ষার মডেলগুলিতে সম্ভাব্য সীমাবদ্ধতার পরামর্শ দেয়। জিন ট্রান্সফরমার, একটি গভীর শিক্ষার কাঠামো যা সিএনএন এবং মনোযোগের প্রক্রিয়াকে একীভূত করে, বিশেষ করে জিন এক্সপ্রেশন বা মেডিকেল ইমেজিং ডেটাসেটের মতো উচ্চ-মাত্রিক ডেটা পরিচালনার ক্ষেত্রে কার্যকর। এটি র্যান্ডম ফরেস্ট মডেলের ব্যাখ্যাযোগ্যতার ব্যবধানকেও সম্বোধন করে, যা গবেষকদের নির্দিষ্ট জিন বা ইমেজিং বৈশিষ্ট্যগুলির মধ্যে অন্তর্দৃষ্টি লাভ করতে দেয় যা শ্রেণিবিন্যাসের সিদ্ধান্তগুলি চালায়। ভবিষ্যতের অধ্যয়নগুলি আরও ব্যাপক ক্যান্সার বিশ্লেষণের জন্য গভীর শিক্ষার মডেলগুলির সাথে ক্লিনিকাল এবং হিস্টোপ্যাথোলজিকাল ডেটা একীভূত করতে পারে।

Conclusion:

The study uses Convolutional Neural Networks (CNNs) and Random Forest models to classify lung cancer subtypes using two datasets: the IQ-OTHNCCD lung cancer dataset and the LC25000 Lung Dataset. The Random Forest model consistently achieved perfect performance, with 100% accuracy, F1-score, precision, and recall. However, CNN's performance on the LC25000 dataset was moderate, indicating potential areas for improvement. The study suggests that traditional machine learning models can be powerful tools in medical imaging analysis.

- ⇒ গবেষণাটি দুটি ডেটাসেট ব্যবহার করে ফুসফুসের ক্যান্সারের উপপ্রকার শ্রেণিবদ্ধ করতে কনভোলুশনাল নিউরাল নেটওয়ার্ক (CNNs) এবং র্যান্ডম ফরেস্ট মডেল ব্যবহার করে: IQ-OTHNCCD ফুসফুসের ক্যান্সার ডেটাসেট এবং LC25000 ফুসফুসের ডেটাসেট। র্যান্ডম ফরেস্ট মডেল ধারাবাহিকভাবে 100% নির্ভুলতা, F1-স্কোর সহ নিখুঁত কর্মক্ষমতা অর্জন করেছে, নির্ভুলতা, এবং প্রত্যাহার। যাইহোক, LC25000 ডেটাসেটে CNN-এর কর্মক্ষমতা মাঝারি ছিল, যা উন্নতির সম্ভাব্য ক্ষেত্রগুলি নির্দেশ করে। গবেষণাটি পরামর্শ দেয় যে ঐতিহ্যগত মেশিন লার্নিং মডেলগুলি মেডিকেল ইমেজিং বিশ্লেষণে শক্তিশালী সরঞ্জাম হতে পারে।

- ⇒ Results Analysis, Discussion, Conclusion, Datasets: Are very important