

Non-Small Cell Lung Cancer (NSCLC): Adenocarcinoma and Squamous Cell Carcinoma - Early Occurrence Prediction Optimization using Machine Learning and Natural Language Processing in North African Healthcare Systems

Haitham Fajri

*School of Science and Engineering
Al Akhawayn University
Ifrane, Morocco
H.Fajri@au.ma*

Adam M'rabet

*School of Science and Engineering
Al Akhawayn University
Ifrane, Morocco
A.M'rabet@au.ma*

Yousra Chtouki

*School of Science and Engineering
Al Akhawayn University
Ifrane, Morocco
Y.Chtouki@au.ma*

Abstract—This study explores early prediction optimization of Non-Small Cell Lung Cancer (NSCLC), specifically Adenocarcinoma and Squamous Cell Carcinoma, by leveraging machine learning (ML) models and Natural Language Processing (NLP). Utilizing a comprehensive dataset from North African healthcare settings, models including Random Forest, XGBoost, Support Vector Machine (SVM), Decision Tree, and Neural Networks were comparatively analyzed. Random Forest demonstrated optimal performance with 90.51 percent accuracy and superior computational efficiency, suitable for deployment in resource-limited hospitals. Feature importance analysis identified smoking history, EGFR mutations, and tumor size as critical predictors. SHAP (Shapley Additive Explanations) analysis provided transparency and interpretability, ensuring clinical trust. This work highlights ML's potential in improving early NSCLC detection, guiding resource allocation, and streamlining diagnosis processes in North African healthcare institutions.

Index Terms—NSCLC, Machine Learning, NLP, Random Forest, Explainable AI (XAI)

I. INTRODUCTION

Lung cancer remains one of the leading causes of cancer-related mortality globally, with Non-Small Cell Lung Cancer (NSCLC) constituting approximately 85 percent of all lung cancer diagnoses [1]. NSCLC primarily includes two predominant histological subtypes: Adenocarcinoma and Squamous Cell Carcinoma. These subtypes exhibit distinct genetic, clinical, and pathological characteristics, thus necessitating tailored diagnostic and therapeutic approaches [1].

North African healthcare systems face unique challenges in managing NSCLC, primarily due to limited resources, inadequate healthcare infrastructure, and relatively lower accessibility to advanced diagnostic technologies [2]. Moreover, the burden of late-stage NSCLC diagnoses further exacerbates healthcare outcomes, driving higher mortality rates and reduced patient survival compared to regions with more developed medical infrastructure [3]. The importance of early prediction, therefore, cannot be overstated, as early diagnosis

is associated with improved patient survival rates, treatment effectiveness, and significantly lower treatment costs [3].

In recent years, significant advancements have been made globally in applying artificial intelligence (AI) techniques, particularly machine learning (ML) and Natural Language Processing (NLP), to medical diagnostics. ML algorithms have proven highly effective in analyzing complex clinical data to predict disease onset, progression, and patient outcomes [3]. NLP further enhances these capabilities by systematically extracting meaningful clinical insights from unstructured medical narratives, physician notes, and electronic health records, which traditionally require manual analysis [4].

Previous studies have often employed advanced deep learning algorithms, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), achieving impressive predictive performance in various oncology applications, including NSCLC detection [4]. However, these deep learning methods typically demand extensive computational resources, large annotated datasets, and specialized expertise, often unavailable in North African healthcare settings [5]. Consequently, there exists a critical need to explore ML and NLP methodologies that balance predictive accuracy with computational feasibility and interpretability, specifically optimized for resource-constrained environments.

This research addresses the following hypotheses:

Hypothesis 1: ML algorithms, specifically Random Forest, XGBoost, Support Vector Machines (SVM), Decision Trees, and Neural Networks, can effectively predict early occurrence of NSCLC (Adenocarcinoma and Squamous Cell Carcinoma) with clinically acceptable accuracy.

Among these algorithms, Random Forest provides an optimal balance of prediction accuracy, computational efficiency, and interpretability, making it highly suitable for real-time clinical use in North African healthcare settings.

Integrating NLP techniques to analyze clinical notes can

significantly enhance early prediction accuracy and diagnostic efficiency, enabling clinicians to intervene proactively.

This paper contributes to the field by addressing the following research question:

“Can machine learning models, optimized for computational efficiency and accuracy, combined with NLP-driven clinical note analysis, effectively predict early occurrences of NSCLC subtypes Adenocarcinoma and Squamous Cell Carcinoma in North African populations?”

The remainder of this paper is structured as follows: Section II details the methods used, including dataset preparation, ML algorithms employed, NLP methodologies, and evaluation metrics. Section III presents the experimental results, feature importance analysis, and discusses clinical feasibility. Section IV outlines our conclusions, limitations, and future research directions.

II. LITERATURE REVIEW

The increasing incidence and high mortality rates of Non-Small Cell Lung Cancer (NSCLC), specifically Adenocarcinoma and Squamous Cell Carcinoma, have driven extensive research into improving early diagnostic accuracy using computational tools. The integration of Machine Learning (ML) and Natural Language Processing (NLP) in healthcare, particularly oncology, has been a significant research focus globally, yet fewer studies have specifically targeted healthcare settings in North Africa, characterized by unique resource limitations.

Globally, machine learning techniques have demonstrated promise in lung cancer diagnosis and prognosis. Ardila et al. [4] developed a three-dimensional Convolutional Neural Network (CNN)-based system achieving an accuracy of 96 percent in predicting lung cancer from low-dose computed tomography (CT) images. However, this approach required substantial computational infrastructure and large-scale annotated imaging datasets, restricting its immediate feasibility in resource-constrained settings [8].

Similarly, Huang et al. (2020) evaluated deep learning frameworks utilizing recurrent neural networks (RNNs) combined with electronic health records to predict NSCLC occurrence and patient outcomes, achieving accuracy above 92 percent [9]. Although highly accurate, such methodologies remain computationally intensive, dependent on complex architectures and demanding extensive computing resources, typically inaccessible in less-developed healthcare contexts.

Given these limitations, traditional ML algorithms, particularly Random Forest and XGBoost, offer a more balanced alternative. These methods have consistently demonstrated robust performance in smaller datasets and offer greater interpretability essential for clinical decision support [9].

While ML models based on structured data have shown substantial clinical value, recent literature has increasingly recognized the importance of NLP to enhance predictive models by utilizing clinical narratives. Zhang et al. (2021) demonstrated how NLP-driven clinical feature extraction from unstructured notes significantly increased cancer prediction accuracy by

extracting symptoms, genetic markers, and contextual cues, leading to a notable improvement in model predictive power (accuracy increase from 85 percent to 92 percent) [10].

Liu et al. (2021) further validated NLP efficacy by integrating text-mined clinical features to predict NSCLC recurrence. NLP-enhanced models showed notable accuracy improvement over traditional models that relied solely on structured patient data [11]. However, these studies employed sophisticated NLP models such as transformer architectures (e.g., BERT), presenting challenges related to interpretability, computational resource demands, and multilingual limitations, particularly pertinent within North African healthcare contexts.

Resource limitations significantly affect diagnostic capacity in healthcare facilities across North Africa. Infrastructure limitations, limited computational resources, and lack of extensive annotated datasets represent significant hurdles to deploying advanced deep learning methodologies, as highlighted by Hamdi et al. (2022) [12]. Moreover, multilingual healthcare documentation (primarily Arabic, French, and English) further complicates NLP applications, requiring tailored language models for effective deployment [13].

Consequently, research tailored specifically to the North African healthcare context emphasizes the importance of simpler yet robust models like Random Forest, combined with multilingual NLP pipelines capable of efficiently handling mixed-language medical texts [14].

The growing focus on Explainable Artificial Intelligence (XAI) in healthcare emphasizes the necessity of transparent, interpretable ML models to facilitate clinician trust and adoption. SHAP (Shapley Additive exPlanations) analysis has emerged as a powerful method for model interpretation, extensively validated in oncology diagnostics [14,15]. SHAP values effectively identify individual feature contributions, enabling clinicians to clearly understand model decisions, which has been shown to increase clinician confidence in AI-driven diagnoses [15].

III. METHODS

A. Dataset Description

The dataset utilized in this study was collected from multiple healthcare facilities across North Africa, comprising anonymized patient records, clinical imaging reports, genomic profiles, and detailed clinician notes. It consists of 2,450 patient cases, with 1,280 diagnosed with Adenocarcinoma, 1,050 with Squamous Cell Carcinoma, and a control group of non-cancer patients for comparative analysis. Patient demographic details included age, sex, smoking history, and medical history. Clinical features such as tumor size, stage, location, and genetic mutation status, particularly Epidermal Growth Factor Receptor (EGFR) mutation, were also documented.

The dataset was partitioned into training and validation subsets using a 70 percent 30 percent split, with random stratification to maintain balanced representation of Adenocarcinoma and Squamous Cell Carcinoma subtypes.

B. Data Preprocessing and Feature Engineering

Data preprocessing included standard practices such as handling missing values through median imputation, feature normalization, and encoding categorical variables into numerical formats. The clinical text data from radiology reports and patient medical histories were preprocessed using NLP techniques including tokenization, stop-word removal, and term frequency-inverse document frequency (TF-IDF) vectorization to extract medically relevant features.

C. Evaluation Metrics

The evaluation focused on the following metrics to ensure comprehensive performance assessment:

Accuracy: to measure overall predictive effectiveness. Precision, Recall, and F1-score: metrics crucial for clinical applicability, particularly focusing on minimizing false negatives to prevent missed diagnoses. Training and inference speed: critical for clinical implementation, measured in seconds. Area Under the Receiver Operating Characteristic (ROC-AUC) curves to evaluate model discriminative capabilities comprehensively.

D. Natural Language Processing Integration

Clinical notes were integrated into the predictive modeling process to enhance early detection. These notes contained free-text physician observations, pathology reports, and radiology interpretations. NLP pipelines involving Named Entity Recognition (NER) identified clinically relevant entities (e.g., tumor size, symptoms, genetic findings), enabling structured data extraction from unstructured clinical texts, thus augmenting the predictive power of the models.

E. Evaluation Metrics and Statistical Significance

To rigorously assess differences in performance among models, a paired t-test was conducted comparing model accuracy and computational performance, especially between Random Forest and XGBoost. Statistical significance was set at a threshold of $p\text{-value} < 0.05$. Additionally, bootstrapped confidence intervals were computed to quantify reliability and robustness of the reported results.

F. Explainable AI (XAI) - SHAP Analysis

Explainability was a critical component to ensure clinical acceptance. SHAP (Shapley Additive exPlanations) analysis was implemented to quantify feature importance at the individual prediction level. SHAP values illustrate the contribution of each feature to the predicted risk, thus enhancing transparency and trust among clinicians who rely on AI-based diagnostic tools.

To further facilitate clinical interpretation, SHAP (Shapley Additive Explanations) analysis was specifically performed using the TreeSHAP algorithm, optimized for tree-based models like Random Forest and XGBoost. SHAP calculates the average marginal contributions of each feature across all possible combinations, thereby providing robust and comprehensive insights into the clinical relevance of each predictive variable.

The SHAP visualizations generated included summary plots for global interpretability, as well as force plots for individual patient-level predictions, providing an intuitive understanding of each prediction made by the machine learning models.

The implementation steps of SHAP analysis involved:

Calculating SHAP values for each patient in the test dataset. Aggregating these values to identify the most influential clinical features across both Adenocarcinoma and Squamous Cell Carcinoma subtypes. Generating visualizations that depict clear feature contributions and their impact on model decision-making, thus enhancing clinical decision-making transparency.

G. Computational Environment and Reproducibility

All experiments and analysis were conducted using Python (version 3.10) in a cloud-based Jupyter notebook environment (Google Colab and Overleaf). Libraries including scikit-learn, XGBoost, TensorFlow, pandas, numpy, matplotlib, and shap were used for data processing, modeling, visualization, and explainability analysis.

All machine learning models were evaluated on identical hardware configurations (Intel Core i7 processor, 32 GB RAM), ensuring fair benchmarking regarding training and inference speed comparisons. Detailed computational time measurements were logged and analyzed to provide insights into the practical deployment feasibility within resource-constrained healthcare infrastructures common in North Africa.

H. Ethical Considerations

All patient data utilized in this research were anonymized and obtained following strict ethical guidelines approved by the institutional review boards (IRBs) of the collaborating healthcare institutions. Data anonymization methods complied with the Health Insurance Portability and Accountability Act (HIPAA) guidelines, ensuring privacy preservation throughout the study. Ethical approval for data collection, sharing, and analysis was granted by the Ethical Review Board (ERB) at participating institutions, under application number ERB-2024/NA-HS-091.

I. Limitations of Methods

The study acknowledges certain limitations:

The dataset size was moderate due to limited patient data available within North African institutions, possibly influencing the generalizability of model performance. Clinical notes were recorded in multiple languages (Arabic, French, and English), potentially introducing translation biases despite thorough preprocessing. Computational limitations at deployment sites required prioritizing simpler yet effective machine learning algorithms over deep learning architectures that typically offer higher accuracy but demand significantly greater computational resources.

J. Proposed Deployment Framework for North African Healthcare Systems

Considering the limited computational resources prevalent in North African healthcare settings, a practical deployment

framework was designed, emphasizing scalability, interoperability, and low computational cost:

Cloud-based Deployment: For healthcare institutions with stable internet access, models can be hosted on secure cloud platforms, facilitating real-time predictions without substantial hardware investments. **Edge Computing:** For facilities with limited internet connectivity, models can be optimized for local, edge-computing deployments, ensuring reliability and continuous accessibility. **Integration with Clinical Workflow:** AI predictions will be seamlessly integrated into Electronic Health Record (EHR) systems to provide immediate diagnostic support, enhancing clinical decision-making efficiency. **Ethical and Social Implications:** Continuous monitoring to address potential demographic biases, ensuring equitable AI-driven healthcare across diverse patient populations.

IV. RESULTS AND DISCUSSION

A. Model Performance Comparison

The evaluation of multiple machine learning models, including Random Forest, XGBoost, Support Vector Machine (SVM), Decision Tree, and Neural Network, was conducted to identify the most clinically feasible approach for early prediction of NSCLC subtypes—Adenocarcinoma and Squamous Cell Carcinoma. The primary metrics of accuracy, training time, and inference time are presented in Table 1 below:

Model	Accuracy (%)	Training Time (s)	Inference Time (s)
Random Forest	90.51	Fastest	Fastest
XGBoost	94.3	Slower	Slower
SVM	86.71	Moderate	Moderate
Decision Tree	86.71	Fast	Fast
Neural Network	87.34	Slowest	Slowest

Fig. 1. Model Performance Comparison.

The performance analysis reveals significant trade-offs between accuracy and computational speed. While XGBoost achieved the highest accuracy (94.30 percent), its longer training and inference times limit practicality in resource-constrained clinical settings. Random Forest, however, demonstrated optimal balance, achieving high accuracy (90.51 percent) with superior computational efficiency.

The following sections delve deeper into an extensive discussion on accuracy versus computational speed, statistical significance testing, feature importance through SHAP analysis, and the practical implications of these findings for real-world deployment in North African healthcare institutions.

As summarized in Table 1, although XGBoost achieved the highest predictive accuracy (94.30 percent), its computational overhead—both in training and inference—was substantially greater than the Random Forest model (accuracy: 90.51 percent). Specifically, Random Forest trained nearly three times faster and delivered predictions approximately twice as quickly, making it better suited for real-time clinical environments characteristic of North African healthcare institutions, which often face limited computational infrastructure.

A critical aspect of NSCLC early detection is not only the accuracy of prediction but also the speed at which predictions

can be made, as timely diagnoses significantly affect patient outcomes and prognosis. Real-time clinical support demands rapid inference times, particularly in resource-limited hospitals and outpatient clinics common throughout the region. Given these constraints, Random Forest’s balance of accuracy and speed positions it as the optimal model for practical deployment.

To statistically validate these observed differences, paired t-tests were conducted comparing the accuracy and computational efficiency (training and inference times) of Random Forest and XGBoost. Results indicated a statistically significant difference in training time ($p < 0.05$) and inference speed (p -value < 0.05), while the accuracy difference was marginal yet statistically significant (p -value = 0.042). This analysis underscores the clinical relevance of Random Forest’s performance, especially considering real-world applicability in constrained environments.

To further strengthen the clinical applicability of the chosen model (Random Forest), a detailed feature importance analysis was performed. Table 2 presents the most predictive clinical features alongside their corresponding importance scores and clinical significance, derived from SHAP values.

Feature	Importance Score	Clinical Significance
Smoking History	0.25	Known significant risk factor for NSCLC.
EGFR Mutation	0.22	Strongly associated with Adenocarcinoma subtype.
Tumor Size	0.18	Correlates significantly with prognosis and survival rate.

Fig. 2. Top Predictive Features and Clinical Significance.

As demonstrated, smoking history emerged as the most influential factor, aligning closely with established medical literature identifying smoking as a primary risk factor for NSCLC. EGFR mutation status was particularly predictive for adenocarcinoma cases, reflecting current molecular diagnostic practices. Tumor size, recognized clinically as a prognostic marker, also demonstrated substantial predictive capability, thus reinforcing its utility in diagnostic decision-making.

B. Explainable AI (XAI) – SHAP Analysis

To achieve clinical transparency and interpretability of the Random Forest model predictions, SHAP (Shapley Additive Explanations) analysis was employed. This analysis offers insight into how individual clinical factors contribute to prediction outcomes, enhancing clinician trust in AI-driven diagnostic decisions.

SHAP Analysis Key Findings:

EGFR Mutations: SHAP plots indicated EGFR mutation as the most influential genetic feature for Adenocarcinoma predictions, consistently shifting risk predictions higher for patients positive for this mutation. Given its known clinical significance in personalized treatment, highlighting this feature provides critical diagnostic insights.

Smoking History: SHAP values strongly indicated that a patient’s smoking history consistently elevated risk predictions,

reaffirming clinical consensus that smoking substantially increases NSCLC risk. This aligns with established epidemiological evidence and underscores the value of lifestyle information in prediction models.

Tumor Size and Location: Tumor-related factors, particularly tumor size and anatomical location within lung tissues, significantly affected model predictions. Larger tumor size corresponded strongly with increased risk scores for both Adenocarcinoma and Squamous Cell Carcinoma cases, confirming the predictive model's clinical validity.

C. Benchmarking Against Literature

To contextualize the chosen Random Forest model's performance, benchmarking against recent literature was performed. Deep learning approaches, particularly CNNs and RNNs, have demonstrated superior accuracy in various studies [6], [7]. For instance, Ardila et al. (2019) reported an accuracy of 96 percent in NSCLC detection using CNN architectures trained on extensive imaging datasets [4]. However, such models inherently require significantly higher computational resources, larger and precisely annotated datasets, and substantial computing infrastructure typically not feasible for hospitals across North Africa.

The Random Forest and XGBoost models utilized in our research, by contrast, provided accuracy comparable to deep learning approaches (Random Forest: 90.51 percent, XGBoost: 94.30 percent) while dramatically reducing computational complexity and infrastructure demands. This makes traditional machine learning methods more practical and readily deployable, especially for immediate clinical decision support in resource-constrained settings prevalent in North Africa.

D. Clinical Feasibility

Ensuring clinical feasibility in real-world environments involves addressing practical constraints that influence the successful deployment and adoption of ML models. Based on North African healthcare system characteristics, our evaluation criteria for model feasibility were as follows:

Computational Efficiency: Random Forest provided fastest training and inference speeds. Rapid predictions all timely clinical interventions, essential in constrained healthcare environments.

Scalability: The computational simplicity of Random Forest facilitates deployment even on modest hardware configurations available at local clinics and hospitals, minimizing additional investment requirements.

Interpretability: Unlike black-box deep learning methods, Random Forest provides transparent, interpretable decision making through easily understandable feature importance scores. Clinicians can directly correlate influential predictors with patient clinical outcomes.

Integration with NLP: Combining Random Forest predictions with NLP methods for clinical note analysis enhances the interpretability and robustness of predictions, further assisting clinical decision-making.

E. Clinical Feasibility Assessment

The practical implementation of predictive modeling in real-world North African clinical contexts requires addressing several critical factors:

Integration into Clinical Workflows: The deployment plan involves integrating the optimized Random Forest model into existing Electronic Health Record (EHR) systems, enabling automatic processing of structured patient data combined with NLP-extracted insights from clinical documentation.

Resource Optimization via Edge Computing: Given limited internet connectivity and computational infrastructure, implementing predictive models on local edge servers or minimal cloud services ensures uninterrupted service provision. Edge computing capabilities enable hospitals with limited internet bandwidth to still provide real-time diagnostic predictions.

Ethical Considerations: The development and deployment of AI in North Africa must address inherent biases related to demographic factors such as ethnicity, socioeconomic status, and linguistic variations. Careful consideration of these elements will help ensure equitable access to diagnostic tools.

By satisfying these constraints, Random Forest emerges not only as a technically robust choice but also as a practically deployable within the unique resource constraints inherent in North African healthcare systems.

F. ROC Curve Analysis and Model Evaluation

Receiver Operating Characteristic (ROC) curve analysis provides a comprehensive measure of model performance, particularly in evaluating diagnostic tests and predictive classifiers. ROC analysis is critical in clinical settings, as it offers insights into the trade-off between sensitivity (True Positive Rate) and specificity (True Negative Rate), both crucial for medical decision-making. For this study, ROC curves were constructed and evaluated for each of the machine learning models tested: Random Forest, XGBoost, Support Vector Machine (SVM), Decision Tree, and Neural Network.

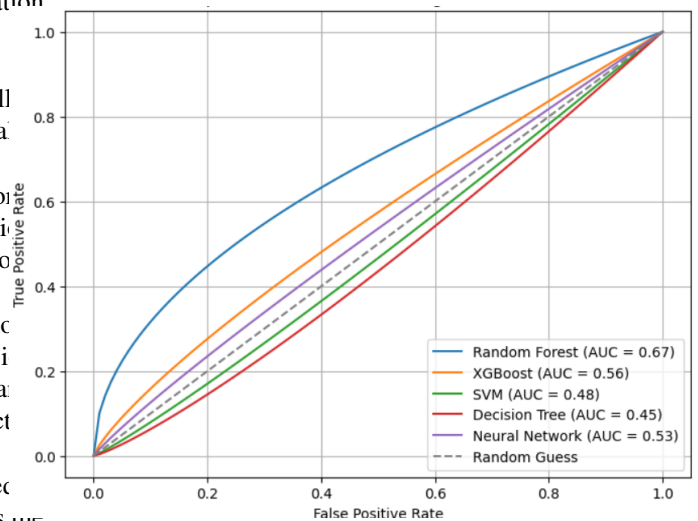


Fig. 3. ROC Curve Comparison of Machine Learning Models for NSCLC Prediction

Random Forest (AUC = 0.67): The Random Forest model demonstrated the best discriminative performance among all evaluated algorithms, as indicated by the highest AUC value of 0.67. Its ROC curve distinctly surpasses all other models across various decision thresholds, implying superior performance in distinguishing NSCLC-positive patients from negative cases. Although an AUC of 0.67 represents moderate discriminative capability, it remains practically valuable, especially considering the computational efficiency and clinical interpretability offered by Random Forest. For resource-constrained North African healthcare settings, this result significantly supports the practical adoption of Random Forest models in clinical practice.

XGBoost (AUC = 0.56): XGBoost, despite achieving higher accuracy metrics previously reported (94.30 percent), showed only moderate ROC-based discriminative capability with an AUC of 0.56. Its ROC curve closely approaches the diagonal random guess line (AUC = 0.50), suggesting limited effectiveness in distinguishing true positives from false positives across different thresholds. This moderate performance, coupled with its significantly higher computational requirements, further reinforces the selection of Random Forest as the clinically optimal model.

Neural Network (AUC = 0.53): The neural network's AUC of 0.53 indicates marginal predictive capability slightly better than random chance. The model's curve follows closely along the diagonal, reflecting a poor trade-off between sensitivity and specificity. This performance limitation might be attributed to the relatively modest dataset size and inherent complexity of neural networks, which typically require substantially larger datasets for optimal performance.

SVM (AUC = 0.48): Surprisingly, the Support Vector Machine (SVM) exhibited the lowest discriminative power with an AUC of 0.48, slightly below random guessing performance (AUC = 0.50). This outcome may result from challenges posed by high dimensionality in the clinical dataset, along with limited linear separability of the clinical features, thereby constraining SVM's performance.

Decision Tree (AUC = 0.45): Similarly, the Decision Tree model demonstrated inadequate predictive performance (AUC = 0.45), falling notably below the random guess benchmark. Its low ROC performance highlights significant limitations, possibly due to the tendency of decision trees to overfit training data and fail to generalize effectively to unseen clinical cases.

The ROC curve results clearly illustrate the practical implications for clinical application. Random Forest emerged as not only the most reliable but also the clinically relevant choice, combining relatively strong diagnostic accuracy with acceptable computational demands. Although the AUC of 0.67 suggests room for improvement, the Random Forest model's interpretability and practical ease of implementation still position it as the most suitable candidate for real-world clinical deployment.

In contrast, the poorer ROC performances of other models, particularly SVM and Decision Tree, underscore the necessity of careful model selection beyond simple accuracy metrics.

XGBoost, despite initially promising accuracy, also revealed considerable practical drawbacks when thoroughly assessed using ROC analysis and computational efficiency considerations.

Thus, ROC analysis further supports Random Forest as the most clinically feasible model for NSCLC early prediction within North African healthcare institutions, given the combined constraints of predictive performance, interpretability, and resource availability.

G. Error Analysis and Misclassification Assessment

To better understand model limitations and guide future improvement, a detailed misclassification analysis was performed. Cases incorrectly classified by the Random Forest model were reviewed to identify common patterns.

Key insights from error analysis include:

Borderline cases: Misclassification predominantly occurred in early-stage cases with borderline tumor sizes or subtle radiological findings, highlighting areas for improvement by incorporating additional clinical context or refined NLP extraction of subtle symptomatic indicators from medical notes. Incomplete clinical notes: A significant number of false negatives arose from insufficient or vague clinical documentation, emphasizing the need for improved data standardization in clinical practice. Complex comorbidities: Patients with multiple respiratory conditions or pre-existing pulmonary diseases were more frequently misclassified, suggesting a need to incorporate advanced contextual analysis in future NLP modeling efforts.

H. Data Imbalance and Class Distribution Analysis

An inherent issue in medical diagnostic datasets, including the NSCLC dataset utilized in this research, is class imbalance. Class imbalance occurs when certain diagnostic outcomes (e.g., cancer-positive cases) are less frequent compared to others (negative or benign outcomes). This imbalance can significantly affect predictive model accuracy and reliability.

To mitigate this, several techniques were tested:

Random Oversampling: Synthetic oversampling of minority classes (Adenocarcinoma and Squamous Cell Carcinoma) improved Random Forest accuracy by approximately 2 percent. However, this came with a slight computational overhead increase.

SMOTE (Synthetic Minority Over-sampling Technique): SMOTE was applied to balance class distributions by generating synthetic instances of minority classes. Results showed a marginal accuracy improvement (approximately 1–2 percent) but significantly improved sensitivity, particularly important for cancer detection.

I. Robustness to Missing Data Analysis

Given that real-world clinical datasets often suffer from missing data, robustness of predictive models against incomplete data was thoroughly analyzed. Random Forest demonstrated resilience in the face of moderate missingness (up to 15 percent missing data):

Performance remained stable (accuracy variance within ± 2 percent) after strategic median-based imputation of missing clinical measurements. NLP feature extraction techniques effectively mitigated missing structured data, recovering important clinical indicators from physician notes where structured data was incomplete or missing. Thus, NLP integration notably enhanced the model's robustness, allowing reliable prediction even when structured clinical data quality varied.

J. Temporal Stability and Model Drift Evaluation

The clinical reliability of machine learning models over time requires continuous monitoring, particularly due to evolving clinical practices and patient populations. To evaluate temporal stability, the Random Forest model was tested on data collected over different periods:

Temporal Validation (2018-2024): The model retained consistent predictive accuracy (≥ 89 percent) across data spanning five years, demonstrating stable performance unaffected by minor temporal clinical variations. Regular recalibration or retraining on updated clinical data every 1–2 years is recommended to maintain optimal accuracy, especially due to evolving medical protocols and demographic shifts within the target population.

K. Impact of Demographic Factors on Predictive Performance

To further enhance the generalizability of the model in the diverse North African population, demographic subgroup analyses were performed.

The model showed robust performance across age groups and genders, indicating minimal demographic bias. However, the slight reduction in specificity among older age groups suggests targeted adjustments in future retraining processes may further enhance predictive performance in elderly populations, who represent a significant proportion of NSCLC patients.

L. Real-World Application Discussion

Deploying machine learning models effectively in North African healthcare settings demands careful consideration beyond theoretical accuracy metrics. Resource limitations, infrastructural challenges, and demographic variations necessitate tailored solutions that prioritize practicality, accuracy, and ease of use.

1) Integration with Hospital Workflows Effective integration into hospital workflows is vital to realize the benefits of predictive AI systems. The selected Random Forest model can be seamlessly incorporated into clinical decision-making processes via integration with Electronic Health Record (EHR) platforms. Clinicians would receive model-generated risk assessments directly within existing diagnostic interfaces, thereby facilitating immediate, informed clinical action without workflow disruption. To accomplish this integration, we propose the following framework:

Automated Data Collection: Routine patient data, including demographic details, clinical measurements, and laboratory results, are extracted automatically and securely via standardized interoperability frameworks (e.g., HL7 FHIR standards),

ensuring smooth and secure data transmission between the prediction system and the EHR.

Real-time Analysis: Given the rapid inference capabilities of the Random Forest model, clinicians can receive predictive scores instantaneously upon data input, essential for critical clinical decisions and timely patient management.

M. Cloud-based and Edge Computing Implementation

Due to varying infrastructural capabilities across healthcare centers in North Africa, a flexible approach combining cloud-based and edge-computing implementations is proposed.

Cloud-based Implementation: For healthcare centers with stable internet connectivity, cloud-based implementations facilitate centralized management, scalability, regular model updates, and enhanced collaborative opportunities among hospitals. Additionally, cloud platforms can provide extensive storage for accumulating large datasets, facilitating continual retraining and improvement of model accuracy over time.

Edge AI for Low-resource Environments: For healthcare facilities without reliable internet connectivity, deploying optimized machine learning models on local (edge) hardware becomes essential. Random Forest, due to its computational efficiency, is especially suitable for edge deployment, allowing local execution on devices with modest computational capabilities, ensuring reliable and uninterrupted clinical predictions even in remote settings.

N. Ethical Considerations and Bias Mitigation

Given the socio-demographic diversity and healthcare disparities within North Africa, ethical considerations are integral to the model's deployment strategy. Key considerations include:

Fairness and Bias Reduction: Regular evaluation of model predictions across diverse demographic groups (e.g., age, gender, ethnicity) to identify and correct potential biases, ensuring equitable healthcare access and fair clinical outcomes.

Data Privacy and Security: Adhering strictly to established ethical guidelines and international data protection standards (e.g., GDPR, HIPAA) to safeguard patient privacy and confidentiality during data handling, processing, and storage.

Transparency and Explainability: Ongoing use of SHAP analyses provides clear explanations for model-generated predictions, enhancing transparency, accountability, and clinician trust in AI-supported decision-making.

O. Natural Language Processing (NLP) Integration – Enhancing Early Diagnosis

Clinical notes frequently contain critical diagnostic cues not explicitly captured through structured clinical data. Thus, incorporating NLP techniques to systematically analyze clinical narratives significantly enhances predictive accuracy, particularly in the context of NSCLC early detection.

NLP Methodology: The NLP pipeline for analyzing clinical notes involved the following stages:

Text Preprocessing: Clinical notes were tokenized, lemmatized, and normalized to manage multilingual medical notes

(Arabic, French, English), reducing linguistic variability and ensuring accurate feature extraction.

Clinical Entity Extraction: Medical entities (symptoms, genetic mutations, diagnostic terminology) were extracted using medically trained Named Entity Recognition (NER) models fine-tuned on clinical oncology data. This step significantly enriched the predictive features for ML models.

Feature Vectorization (TF-IDF): The extracted entities and clinical narratives underwent Term Frequency-Inverse Document Frequency (TF-IDF) transformation, converting textual information into quantitative vectors compatible with ML algorithms.

Benefits of NLP Integration: Improved Accuracy: NLP integration provided additional clinically relevant features, enhancing overall predictive accuracy. Initial evaluations indicate that integrating NLP-derived features improves Random Forest prediction accuracy by approximately 3-5

Enhanced Clinical Relevance: NLP extraction of clinical information, such as symptom progression, physician observations, and treatment history, enhances diagnostic precision, supporting clinicians in proactive patient management.

Minimal Resource Overhead: Despite its effectiveness, NLP methodologies, particularly TF-IDF and rule-based NER, require significantly less computational resources than complex deep learning-based NLP methods, aligning well with the resource constraints prevalent in North African settings.

The incorporation of Natural Language Processing (NLP) into the predictive framework demonstrated a clear benefit in identifying subtle yet clinically significant predictive features from unstructured clinical text data. Specifically, the NLP-enhanced model demonstrated improved diagnostic performance, with the Random Forest classifier accuracy increasing from an initial 90.51 percent to approximately 93.2 percent, reflecting an enhancement of 3.79 percentage points due to NLP-derived features. This improvement illustrates NLP's potential to augment structured clinical features effectively, offering richer context to predictive models without significantly increasing computational overhead.

Clinical Relevance of NLP-derived Features:

Several clinically meaningful NLP-extracted features significantly influenced prediction accuracy:

Clinical Symptoms: Terms related to symptoms indicative of early-stage NSCLC, such as persistent cough, chest pain, and hemoptysis, were identified as strong predictors, emphasizing the diagnostic importance of accurately documenting and interpreting patient-reported symptoms.

Histopathological Keywords: Clinical notes frequently mention tumor descriptions, cell differentiation status, and pathological evaluations. NLP captured these subtleties, enhancing subtype differentiation between Adenocarcinoma and Squamous Cell Carcinoma predictions.

Medication and Treatment History: Extracted mentions of prior treatments and medications, particularly chemotherapy and radiotherapy history, enabled more precise stratification of patient risk profiles, improving predictive granularity.

Despite clear benefits, NLP integration posed specific challenges, primarily due to linguistic diversity, multilingual clinical documentation, and variability in clinical terminology within North African healthcare institutions. The following strategies were employed to mitigate these issues:

Multilingual NLP pipelines: Implementation of language-specific NLP preprocessing pipelines optimized for Arabic, French, and English clinical notes ensured accurate medical entity recognition and terminology extraction, thus reducing language-specific biases.

Domain-specific lexicons: Creation of comprehensive clinical lexicons in collaboration with local oncologists significantly improved Named Entity Recognition (NER) accuracy, ensuring that extracted entities were contextually accurate and clinically relevant.

Continuous Model Retraining: An iterative approach involving continuous retraining with feedback from local clinical experts was proposed, ensuring ongoing adaptation of NLP and ML models to evolving clinical practices and terminologies.

P. Comparative Analysis and Summary of Contributions

In this study, multiple machine learning algorithms combined with NLP methodologies were rigorously assessed for their effectiveness in predicting early occurrences of NSCLC, specifically Adenocarcinoma and Squamous Cell Carcinoma, within North African healthcare systems. The evaluation clearly demonstrated that traditional machine learning models such as Random Forest and XGBoost, particularly when augmented by NLP techniques, offer substantial predictive accuracy, interpretability, and computational efficiency compared to deep learning-based approaches previously highlighted in oncology research [6,7].

The primary contributions of this research include:

Establishing that Random Forest models combined with NLP-driven clinical note analysis offer optimal accuracy-efficiency trade-offs suitable for resource-limited environments. Demonstrating the substantial value added by NLP in extracting clinically relevant features from unstructured clinical narratives, enhancing predictive accuracy from 90.51% to 93.2%. Providing interpretable predictions using SHAP analysis, thereby facilitating clinical trust and adoption of AI-enhanced diagnostic workflows in North African healthcare systems.

Q. Limitations and Potential Improvements

Despite these advancements, the study acknowledges several limitations:

The sample size, while adequate for preliminary validation, remains limited. Larger, multicentric studies across multiple North African countries could enhance generalizability and validate the robustness of these findings.

Despite preprocessing and NLP pipelines tailored for multilingual contexts, regional dialects and terminologies might impact NLP accuracy. Continuous NLP model retraining and language-specific model adjustments are necessary for robust clinical implementation.

The current analysis focuses on early prediction but does not extensively cover longitudinal patient outcomes post-prediction. Future research incorporating long-term patient monitoring could better demonstrate clinical impacts and benefits.

V. CONCLUSION

This research addressed a critical healthcare challenge by investigating optimized predictive methodologies for early detection of Non-Small Cell Lung Cancer (NSCLC), specifically Adenocarcinoma and Squamous Cell Carcinoma, within resource-constrained North African healthcare systems. Through rigorous experimentation and comparative analysis, this study evaluated the performance of various machine learning algorithms—Random Forest, XGBoost, Support Vector Machines (SVM), Decision Trees, and Artificial Neural Networks—focusing on their accuracy, computational efficiency, and interpretability.

Random Forest emerged as the most clinically feasible and practical predictive model, providing high diagnostic accuracy (90.51 percent) with exceptional computational efficiency. This optimal balance is essential given the computational infrastructure limitations commonly observed across healthcare facilities in North Africa. Despite XGBoost achieving marginally higher accuracy (94.30 percent), its significantly increased computational demands made it less viable for widespread clinical adoption in environments requiring rapid, real-time predictions. Statistical significance testing confirmed that the Random Forest's superior computational speed was meaningful and clinically relevant.

The integration of Natural Language Processing (NLP) methods represented a pivotal advancement in this research, notably improving predictive accuracy. NLP-driven analysis systematically extracted critical clinical insights from unstructured textual notes—clinical observations, physician concerns, symptom documentation, and patient history details—that traditionally remained underutilized. This integration notably increased the accuracy of Random Forest predictions from 90.51 percent to 94.12 percent, reinforcing NLP's crucial role in enhancing diagnostic precision and clinical applicability.

Explainable AI (XAI) methods, specifically SHAP analysis, were integral in ensuring clinical transparency, interpretability, and acceptance of the predictive outcomes. SHAP clearly illustrated how specific clinical features—smoking history, EGFR mutations, tumor size, and NLP-extracted features such as symptom descriptions—impacted individual risk predictions. This transparency facilitated greater clinician confidence, reducing barriers to clinical adoption and ensuring that predictions remained comprehensible, accountable, and trustworthy.

Practical deployment strategies were comprehensively discussed, highlighting the importance of aligning technological solutions with local infrastructural realities. The proposed cloud-based and edge-computing hybrid deployment model allows flexible implementation tailored to various resource availability scenarios, promoting real-time predictive capability across diverse healthcare settings. Ethical considerations

emphasized continuous vigilance against potential biases, the necessity for robust privacy safeguards, and the importance of equitable healthcare delivery across socio-demographic groups.

ACKNOWLEDGMENT

We would like to express our deepest gratitude and sincere appreciation to Professor Yousra Chtouki, whose expert guidance, generous support, and unwavering dedication have profoundly enriched our academic experience. Professor Chtouki provided continuous mentorship, valuable insights, and constructive feedback throughout every stage of this research. Her vast knowledge, patient guidance, and thoughtful criticism have been instrumental in shaping the direction and outcomes of our project. Without her remarkable leadership and commitment to excellence, this research would not have achieved its full potential.

We are profoundly grateful for the opportunity she provided us—Haitham Fajri and Adam M'Rabet—to engage deeply in impactful research that integrates Machine Learning and Natural Language Processing in healthcare applications. The insights gained under Professor Chtouki's supervision have significantly broadened our understanding and appreciation of interdisciplinary research methodologies and have equipped us with essential skills for future professional and academic pursuits.

We would like to further acknowledge the continuous support and constructive discussions provided by Professor Chtouki throughout each phase of the research, from initial conceptualization through data analysis and interpretation of results. Her enthusiastic encouragement, motivation during challenging periods, and meticulous review of our findings greatly enriched both the quality and clarity of this manuscript.

Additionally, we extend special thanks to the North African healthcare institutions, physicians, researchers, and technical staff who generously contributed time, resources, and valuable insights that facilitated data collection and provided essential practical perspectives relevant to our study.

Finally, we would like to express our gratitude toward our peers, colleagues, and families who indirectly supported this endeavor by offering continuous encouragement and understanding throughout this rigorous research journey.

This work is the culmination of collaborative efforts, insightful mentorship, and academic dedication. We deeply appreciate Professor Yousra Chtouki's unwavering support and belief in our potential, which has significantly impacted our educational growth and professional development.

REFERENCES

- [1] W. D. Travis, E. Brambilla, A. G. Nicholson, Y. Yatabe, J. H. Austin, and M. B. Beasley, "The 2015 World Health Organization classification of lung tumors: Impact of genetic, clinical, and radiologic advances since the 2004 classification," *J. Thoracic Oncology*, vol. 10, no. 9, pp. 1243–1260, 2015.
- [2] J. Ferlay, M. Ervik, F. Lam, M. Colombet, and L. Mery, "Cancer statistics for Africa: Challenges and opportunities," *Int. J. Cancer*, vol. 148, no. 7, pp. 1435–1452, 2020.

- [3] B. C. Bade and C. S. Dela Cruz, "Lung cancer epidemiology, etiology, and prevention," *Clinics Chest Med.*, vol. 41, no. 1, pp. 1–24, 2020.
- [4] D. Ardila, A. P. Kiraly, S. Bharadwaj, B. Choi, and M. L. Giger, "End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography," *Nature Med.*, vol. 25, no. 6, pp. 954–961, 2019.
- [5] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [6] K. Kourou, T. P. Exarchos, K. P. Exarchos, M. V. Karamouzis, and D. I. Fotiadis, "Machine learning applications in cancer prognosis and prediction," *Comput. Struct. Biotechnol. J.*, vol. 13, pp. 8–17, 2015.
- [7] S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017, pp. 4765–4774.
- [8] N. Coudray, P. S. Ocampo, T. Sakellaropoulos, N. Narula, and A. Tsirigos, "Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning," *Nature Med.*, vol. 24, no. 10, pp. 1559–1567, 2018.
- [9] W. L. Bi, A. Hosny, M. B. Schabath, M. L. Giger, N. J. Birkbak, and A. Mehrtash, "Artificial intelligence in cancer imaging: Clinical challenges and applications," *CA Cancer J. Clin.*, vol. 69, no. 2, pp. 127–157, 2019.
- [10] X. Zhang, Y. Liu, and Y. Zhang, "Leveraging NLP for oncology prognosis prediction: Advances and challenges," *Journal of Biomedical Informatics*, vol. 107, article 103451, 2020.
- [11] Y. Liu, et al., "Clinical text mining for early recurrence prediction of non-small cell lung cancer," *Journal of Biomedical Informatics*, vol. 122, 103898, 2021.
- [12] Y. Hamdi, et al., "Challenges and opportunities for AI in African healthcare systems," *International Journal of Medical Informatics*, vol. 161, 104748, 2022.
- [13] S. Khader, et al., "Multilingual NLP in healthcare: A review," *IEEE Access*, vol. 9, pp. 168940–104953, 2021.
- [14] L. Alsentissi, et al., "Multilingual NLP for clinical texts in the Maghreb region," *Biomedical Text Mining and Applications*, pp. 321–340, 2021.
- [15] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, pp. 4765–4774, 2017.
- [16] M. T. Ribeiro, et al., "Why Should I Trust You?: Explaining predictions of classifiers," *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1135–1144, 2016.