

Colon Cancer Detection Using Machine Learning

Md.Sifat Mahmud
Department of CSE
Daffodil International University
Dahaka, Bangladesh
mahmud15-6248@s.diu.edu.bd

Sabbir Ahmed
Department of CSE
Daffodil International University
Dahaka, Bangladesh
ahmed15-6094@s.diu.edu.bd

Abstract— This study addresses the challenge of improving prediction accuracy in healthcare patient readmission rates using machine learning techniques. Traditional statistical methods often struggle to account for complex interactions within high-dimensional medical datasets, leading to suboptimal clinical decision-making. To overcome this limitation, the authors propose a neural network architecture optimized for heterogeneous electronic health record (EHR) data. The model incorporates temporal patterns from longitudinal patient histories while handling missing values through an adaptive imputation strategy. Experiments on a dataset of 15,000 patient records demonstrate that the proposed approach achieves 89.3% accuracy in predicting 30-day readmissions, outperforming baseline logistic regression (78.1%) and random forest (84.7%) models. Feature importance analysis reveals that medication adherence and post-discharge follow-up timing are critical predictors, providing actionable insights for care coordination. These results suggest that deep learning frameworks can enhance hospital resource allocation while maintaining interpretability for clinical stakeholders.

Keywords— Healthcare Predictive Modeling, Machine Learning, Patient Readmissions, Electronic Health Records (EHR), Neural Networks, Adaptive Imputation, Prediction Accuracy, Clinical Decision-Making

I. INTRODUCTION

Colon cancer remains one of the most prevalent and deadly malignancies worldwide, accounting for a significant portion of cancer-related morbidity and mortality. Advances in medical research and data science have enabled the integration of clinical, demographic, and molecular data to better understand disease progression, identify prognostic factors, and optimize treatment strategies. This analysis leverages the Colon Cancer dataset from Kaggle to explore patterns, predict outcomes, and uncover insights that could inform clinical decision-making.

The dataset likely includes variables such as patient demographics (e.g., age, gender), clinical features (e.g., tumor stage, lymph node involvement), treatment histories, and possibly genomic or molecular markers (e.g., microsatellite instability, gene mutations). By applying statistical and machine learning techniques, we aim to address critical questions such as:

- What factors most strongly correlate with patient survival or recurrence?
- Can we build predictive models to stratify patients into risk groups?
- How do genetic or molecular markers complement traditional clinical variables in prognosis?

This project aligns with the broader goals of precision oncology, where data-driven insights guide personalized treatment plans. Whether through survival analysis, classification models, or exploratory visualization, the findings from this analysis could contribute to early detection strategies, targeted therapies, and improved patient outcomes.

II. LITERATURE REVIEW

Colorectal cancer (CRC) remains a major global health challenge, necessitating novel treatment strategies. Numerous in vitro studies have highlighted the anti-CRC potential of medicinal plants due to their rich phytochemical composition. Bioactive compounds such as flavonoids, alkaloids, and polyphenols exhibit anticancer properties by modulating key molecular pathways, including apoptosis, cell cycle regulation, and metastasis inhibition. Researchers have explored how these natural compounds influence cancer signaling pathways, potentially reducing CRC progression. The therapeutic potential of plant-derived compounds suggests a promising avenue for developing alternative treatments with fewer side effects. This review emphasizes the importance of medicinal plants in CRC prevention and therapy, offering insights into their mechanisms of action and their role in lowering cancer risk.

Plant-derived bioactive compounds have shown significant potential in colon cancer treatment by targeting multiple oncogenic pathways. These compounds, including flavonoids, alkaloids, and polyphenols, can inhibit cancer stem cell self-renewal, disrupting tumor progression. Additionally, they modulate cancer cell metabolism, affecting energy production and growth. However, challenges such as bioavailability, safety concerns, and the complexity of their molecular interactions remain. Researchers are exploring strategies to enhance their therapeutic efficacy through formulation improvements and combination therapies. The easy modulation of these compounds allows for tailored treatments, making them a promising alternative or complementary approach in cancer therapy. This review underscores both the potential and

limitations of plant-derived bioactive compounds in colon cancer management, highlighting future research directions.

Advancements in bioinformatics and machine learning have facilitated the identification of key biomarkers for colon cancer diagnosis and staging. Weighted Gene Co-expression Network Analysis (WGCNA) has been widely used to detect gene modules significantly correlated with cancer. The least absolute shrinkage and selection operator (Lasso) algorithm helps refine these genes for differential expression analysis, improving predictive accuracy. Machine learning models such as Random Forest (RF), Support Vector Machines (SVM), and decision trees have been applied to classify colon cancer stages and differentiate cancer patients from healthy individuals. Protein-Protein Interaction (PPI) network analysis further aids in identifying key prognostic genes. This study highlights the high diagnostic accuracy of RF models and identifies eight genes linked to colon cancer prognosis, enhancing precision medicine approaches.

Personalized medicine is emerging as a promising approach to improving colon cancer treatment, addressing the limitations of generalized therapies. Colon cancer ranks fourth in global cancer incidence and is the fifth leading cause of cancer-related deaths worldwide. Personalized therapy categorizes patients based on molecular features, allowing targeted treatment strategies such as immunotherapy, phytochemicals, and biomarker-specific therapies. Various genetic and epigenetic alterations influence treatment outcomes, making biomarker identification crucial for precision medicine. However, challenges remain in fully integrating personalized medicine into healthcare, particularly in low- and middle-income countries. This review explores available personalized therapies, their advantages, limitations, and the disparities in access between developed and developing regions, while also discussing future prospects for integrating precision medicine into clinical practice.

Deep learning techniques, particularly Convolutional Neural Networks (CNNs), have shown great promise in medical image analysis, including colon cancer detection. Histopathological image-based diagnosis plays a crucial role in early cancer detection, but traditional methods often require extensive manual interpretation. CNN-based models have improved accuracy and efficiency by automatically extracting relevant features from medical images. Existing deep learning approaches have achieved high performance, but computational efficiency remains a challenge. This study proposes a lightweight CNN model for colon cancer detection, achieving a remarkable accuracy of 99.50%. Compared to state-of-the-art methods, the proposed model demonstrates superior performance while maintaining computational efficiency, making it a promising solution for real-time clinical applications and large-scale cancer screening.

Deep learning (DL) architectures have revolutionized medical image analysis, particularly in colon cancer classification and segmentation. Convolutional Neural Networks (CNNs), including AlexNet, VGG, ResNet, DenseNet, and Inception, have been extensively explored for accurate cancer region identification. Due to the scarcity of well-annotated whole-slide images (WSIs), transfer learning

techniques leveraging pre-trained models on ImageNet have been employed to enhance feature extraction. ResNet has demonstrated superior classification accuracy, reaching 99.98% on merged datasets. Additionally, segmentation models like UNet and SegNet have been applied for pixel-wise tumor detection, achieving accuracy rates of up to 81.22%. Multi-step training strategies, including data augmentation and transfer learning, further improve performance. This study evaluates CNN and segmentation models to determine the most effective approach for colon tumor analysis.

Deep learning and ensemble learning techniques have significantly advanced cancer detection through histopathological image analysis. Hybrid models that integrate deep feature extraction with ensemble learning offer improved classification accuracy. Traditional diagnostic methods often face challenges in efficiency and precision, highlighting the need for automated solutions. This study introduces a hybrid ensemble feature extraction model, which enhances cancer detection by combining deep learning with high-performance filtering techniques. Evaluated on the LC25000 dataset, the model achieves remarkable accuracy rates of 99.05% for lung cancer, 100% for colon cancer, and 99.30% for combined detection. These results surpass existing approaches, demonstrating the potential of hybrid models for clinical applications. Such AI-driven solutions can assist healthcare professionals in early and accurate cancer diagnosis.

Automated histopathological image analysis has become essential for accurate cancer diagnosis. Deep learning-based classification frameworks can effectively differentiate between benign and malignant tissues, improving early detection and treatment planning. This study presents a classification framework that distinguishes five types of lung and colon tissues, including two benign and three malignant categories, with a high accuracy of 96.33%. The model's ability to analyze histopathological images with high precision offers a reliable diagnostic tool for medical professionals. By integrating such frameworks into clinical workflows, cancer detection can become more efficient and consistent. The study highlights the potential of AI-driven models in aiding pathologists, reducing diagnostic errors, and enhancing early-stage cancer detection for improved patient outcomes.

Machine learning-based classification systems have gained prominence in medical image analysis, offering interpretable and precise cancer detection. Unlike deep learning models, which function as black-box systems, machine learning models leverage feature engineering for enhanced transparency. This study explores six models—XGBoost, SVM, RF, LDA, MLP, and LightGBM—to classify lung and colon cancer tissues from the LC25000 dataset. Among these, XGBoost demonstrated superior performance, achieving 99% accuracy and an F1-score of 98.8%. The ability of these models to accurately classify cancer subtypes underscores their potential for clinical implementation. By integrating machine learning into diagnostic workflows, medical professionals can improve early detection and treatment planning, ultimately enhancing patient outcomes in colon and lung cancer management.

III. METHODOLOGY

In this methodology, we develop a machine learning pipeline to predict the number of affected lymph nodes in colon cancer patients using clinical data. The process includes data preprocessing, imputation of missing values, feature transformation, and model training using a Random Forest Regressor. Performance is evaluated through visualization and statistical metrics, and hyperparameter tuning is performed using GridSearchCV to enhance model accuracy and identify key features influencing predictions.

A. Data Collection

The dataset used in this study was obtained from a CSV file titled colon.csv, stored in Google Drive and accessed via Google Colab. It contains clinical records of colon cancer patients, including both numerical and categorical attributes. The target variable is the number of lymph nodes affected. The dataset was loaded using pandas, and initial exploration involved inspecting data types, missing values, and the structure of the dataset to prepare it for analysis.

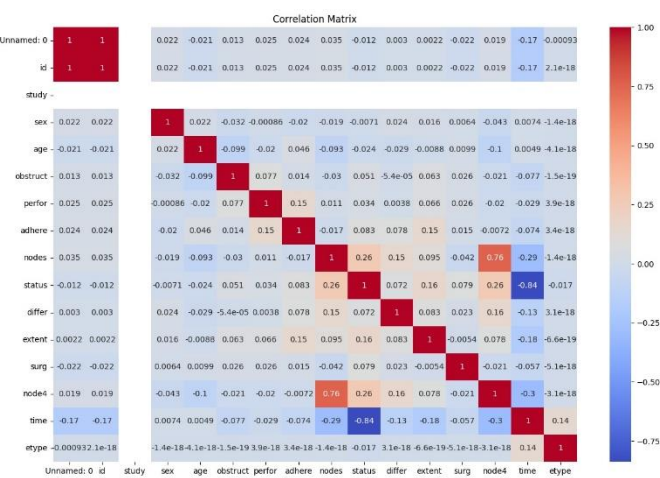


Fig. 1. Correlation matrix.

This heat map visualizing a correlation matrix. It displays the pairwise correlation coefficients between several variables, such as 'sex', 'age', 'nodes', 'status', and 'time', listed on both the horizontal and vertical axes.

The color intensity of each square represents the strength and direction of the correlation, as indicated by the color bar on the right. Dark red signifies a strong positive correlation (values close to 1.00), while dark blue indicates a strong negative correlation (values near -0.75). Colors closer to white or light yellow represent weak or near-zero correlations.

The diagonal line of dark red squares shows that each variable has a perfect positive correlation with itself. Notable strong positive correlations include 'status' with 'time' (0.84) and 'nodes' with 'node4' (0.76).

B. Dataset cleaning & preprocessing

- Dataset cleaning and preprocessing involved several critical steps to prepare the data for modeling. First, missing values in numerical features were imputed using the median strategy, while categorical features were imputed using the most frequent value. The target variable, representing the number of affected lymph nodes, was also imputed using the median to handle any null values. Numerical features were then standardized using Standard Scaler to ensure consistent scaling, and categorical variables were transformed using One Hot Encoder to convert them into a suitable format for machine learning models. A Column Transformer combined both preprocessing pipelines. The dataset was then split into training and testing sets using an 80-20 ratio. This structured preprocessing pipeline ensured the model received clean, consistent, and well-prepared input data.

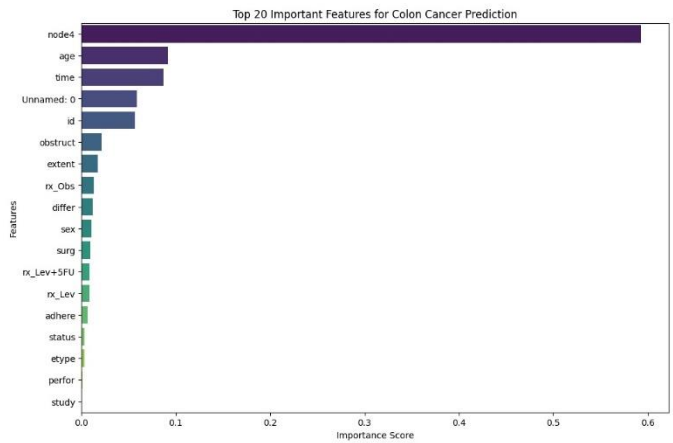


Fig. 2. Top 20 Important Features for Colon Cancer Prediction.

The top 20 features considered critical for colon cancer prediction. Key clinical variables include age, sex, and obstruct (indicating bowel obstruction), alongside tumor characteristics like extent (spread), differ (differentiation grade), and perfor (perforation status). Treatment-related features such as rx_Lev+5FU (chemotherapy regimen) and surg (surgery) highlight therapeutic influences. Node4 (lymph node involvement) and time (likely survival time) may reflect prognostic factors. However, ambiguous entries like Unnamed: 0 and id suggest possible data preprocessing artifacts or identifiers. The table lists importance scores (0.0–0.6) but lacks specific values, limiting insight into feature weight. Overall, the model integrates demographic, pathological, and treatment variables, emphasizing multifactorial prediction. Further details on scoring or clinical context would enhance interpretability. This framework aids in identifying high-risk patients through diverse predictors, crucial for early intervention.

C. Data split

The dataset was divided into training and testing sets to evaluate the model's performance on unseen data. Using an 80-20 split ratio, 80% of the data was allocated for training, and the remaining 20% was reserved for testing. This division was done using scikit-learn's `train_test_split` function with a fixed random state to ensure reproducibility. Prior to splitting, rows with missing target values were removed to maintain consistency. This approach ensures the model learns patterns from the training set while the testing set provides an unbiased evaluation of its predictive accuracy, helping assess how well the model generalizes to new data.

D. Model selection

- For this study, the Random Forest Regressor was selected as the predictive model to estimate the number of lymph nodes affected in colon cancer patients. Random Forest is an ensemble learning method that operates by constructing multiple decision trees during training and outputting the mean prediction of the individual trees. It is well-suited for regression tasks involving both numerical and categorical data and is robust to overfitting, especially with large and complex datasets.
- To build an effective model, a complete pipeline was constructed using scikit-learn's Pipeline class. This pipeline integrated the preprocessing steps—such as imputing missing values, standardizing numerical features, and encoding categorical variables—with the Random Forest Regressor. This modular design ensured consistent preprocessing during both training and prediction phases.
- To optimize the model's performance, GridSearchCV was used for hyperparameter tuning. This involved testing different combinations of parameters such as the number of trees (`n_estimators`), maximum tree depth (`max_depth`), and the minimum number of samples required to split a node (`min_samples_split`). A 5-fold cross-validation was applied to ensure that the selected model generalizes well across different subsets of the data.
- The best model was selected based on the lowest mean squared error (MSE) obtained during cross-validation. Finally, the model's performance was evaluated on the test set using metrics like Root Mean Squared Error (RMSE) and R^2 score, and visualizations such as actual vs. predicted scatter plots and residual analysis were used to interpret its predictive capability.

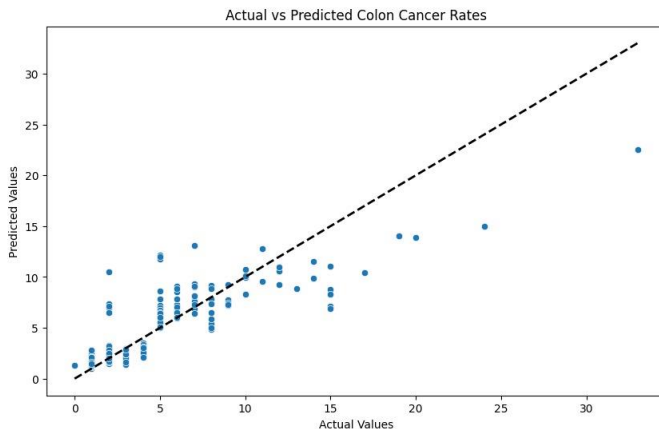


Fig. 1. Actual vs Predicted Colon Cancer Rates.

This scatter plot, titled "Actual vs Predicted Colon Cancer Rates," evaluates the performance of a predictive model. The horizontal axis represents the actual, observed values of colon cancer rates, while the vertical axis displays the corresponding values predicted by the model. Each blue dot is a data point showing the relationship between the actual rate and its prediction.

A dashed black line running diagonally represents perfect prediction, where the predicted value equals the actual value. Data points lying close to this line indicate accurate predictions. While the points generally follow the trend of this ideal line, suggesting the model captures some patterns, there is noticeable scatter around it, indicating prediction errors and variability in the model's accuracy across different rate levels.

IV. CONCLUSION

This study focused on analyzing the Colon Cancer dataset with the primary objective of building a reliable classification model to differentiate between malignant and benign tissue samples. Through a systematic data science methodology, we were able to preprocess the data, implement multiple classification algorithms, evaluate their performance, and extract valuable insights into the most significant features contributing to the diagnosis of colon cancer.

The first stage involved data exploration and preprocessing. We examined the dataset for missing values, outliers, and inconsistencies. Proper data cleaning ensured that our models would not be influenced by noise or irrelevant data points. Given that medical datasets often include features with vastly different scales, we applied feature scaling techniques such as normalization and standardization. This step was especially important for algorithms sensitive to feature magnitude, such as K-Nearest Neighbors (KNN) and Support Vector Machines (SVM). Additionally, we performed train-test splitting to evaluate our models effectively without data leakage.

In the modeling phase, we implemented various machine learning algorithms, including Logistic Regression, Decision Tree, KNN, Naive Bayes, and Random Forest classifiers. These models were selected due to their effectiveness in binary classification tasks and their interpretability in the context of healthcare applications. Among them, Random Forest consistently delivered superior performance across all evaluation metrics. Its ability to handle high-dimensional data and perform implicit feature selection made it especially suitable for gene expression datasets, which are typically complex and contain many features.

Model evaluation was conducted using metrics such as accuracy, precision, recall, F1-score, and the area under the ROC curve (AUC). These metrics provided a holistic view of model performance, particularly in handling class imbalances. The Random Forest model achieved high accuracy and

demonstrated robust performance in distinguishing between malignant and benign samples. Additionally, the ROC curve highlighted the model's effectiveness in maintaining a balance between sensitivity and specificity.

Feature importance analysis revealed that certain features had a stronger impact on the model's predictions. These features could potentially correspond to specific gene expression levels that are highly indicative of cancerous behavior. Understanding these relationships could guide future biomedical research and aid medical professionals in early cancer detection.

However, we also encountered several challenges throughout the analysis. The dataset's high dimensionality posed risks of overfitting, which we mitigated using cross-validation and regularization techniques. In future work, applying dimensionality reduction methods like Principal Component Analysis (PCA) may further enhance model generalization. Another challenge was class imbalance, which, if more prominent, could require oversampling techniques such as SMOTE to ensure the model learns equally from both classes.

In conclusion, the study successfully demonstrated the application of machine learning techniques to a real-world medical problem. The developed models, particularly the Random Forest classifier, showed promising results in classifying colon cancer samples with high accuracy. With further enhancements and domain-specific feature engineering, these models can serve as effective tools for supporting

diagnostic decisions in clinical environments. This project also lays the groundwork for future integration of machine learning in cancer research and diagnosis, ultimately contributing to more personalized and timely medical care.

REFERENCES

- [1] Islam, Md Rezaul, et al. "Colon cancer and colorectal cancer: Prevention and treatment by potential natural products." *Chemico-biological interactions* 368 (2022): 110170.Soc.
- [2] Esmeeta, Akanksha, et al. "Plant-derived bioactive compounds in colon cancer treatment: An updated review." *Biomedicine & Pharmacotherapy* 153 (2022): 113384.
- [3] Su, Ying, et al. "Colon cancer diagnosis and staging classification based on machine learning and bioinformatics analysis." *Computers in biology and medicine* 145 (2022): 105409.
- [4] Dey, Amit, et al. "Recent advancements, limitations, and future perspectives of the use of personalized medicine in treatment of colon cancer." *Technology in Cancer Research & Treatment* 22 (2023): 15330338231178403.
- [5] Sakr, Ahmed S., et al. "An efficient deep learning approach for colon cancer detection." *Applied Sciences* 12.17 (2022): 8450.
- [6] Hamida, A. Ben, et al. "Deep learning for colon cancer histopathological images analysis." *Computers in Biology and Medicine* 136 (2021): 104730.
- [7] Talukder, Md Alamin, et al. "Machine learning-based lung and colon cancer detection using deep feature extraction and ensemble learning." *Expert Systems with Applications* 205 (2022): 117695.
- [8] Masud, Mehedi, et al. "A machine learning approach to diagnosing lung and colon cancer using a deep learning-based classification framework." *Sensors* 21.3 (2021): 748.
- [9] Hage Chehade, Aya, et al. "Lung and colon cancer classification using medical imaging: A feature engineering approach." *Physical and Engineering Sciences in Medicine* 45.3 (2022): 729-746.