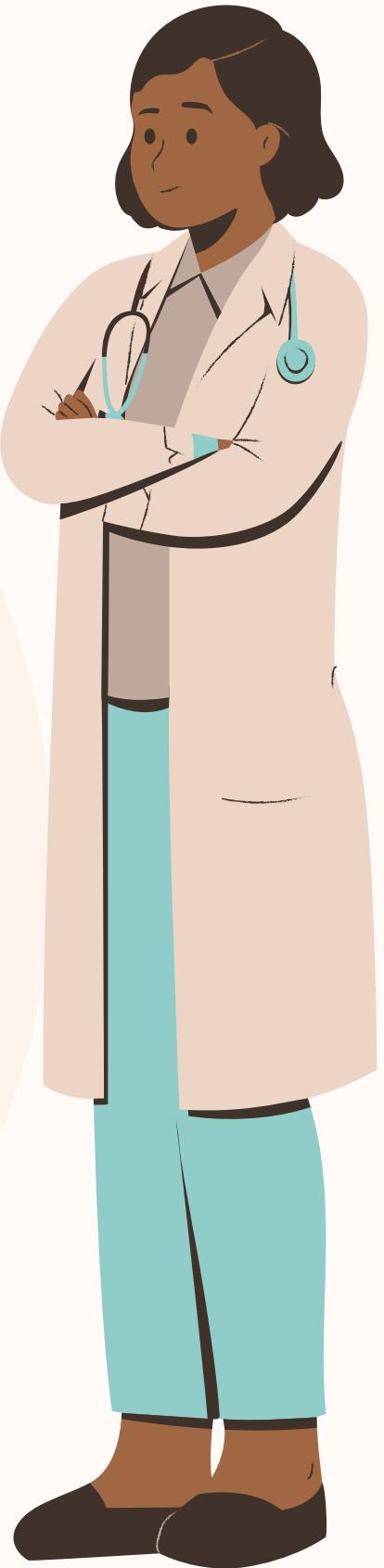


LUNG CANCER SURVIVAL ANALYSIS USING IBM SPSS MODELER



TEAM MEMBERS

- **Anjali Sharma 220239**
 - **Diksha Khangarot 220244**
- Kiran 220247**



INTRODUCTION TO SURVIVAL ANALYSIS

**An exploration of
survival analyses
for cancer**

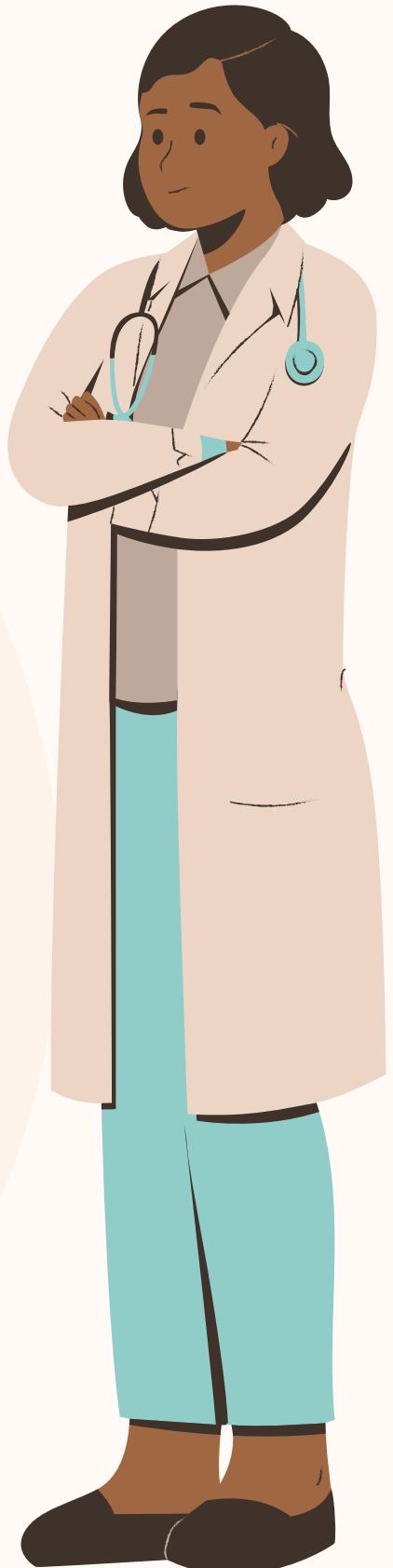
Survival analysis is a statistical method used to estimate the time until an event occurs, like death or relapse, especially in medical research. In this project, we apply survival analysis to lung cancer, one of the deadliest forms of cancer, to understand which factors influence patient survival. The goal is to use IBM SPSS Modeler to analyze a dataset of lung cancer patients and identify key variables such as treatment type, tumor type, and patient condition that impact survival outcomes. This can help in making better clinical decisions and improving patient care.



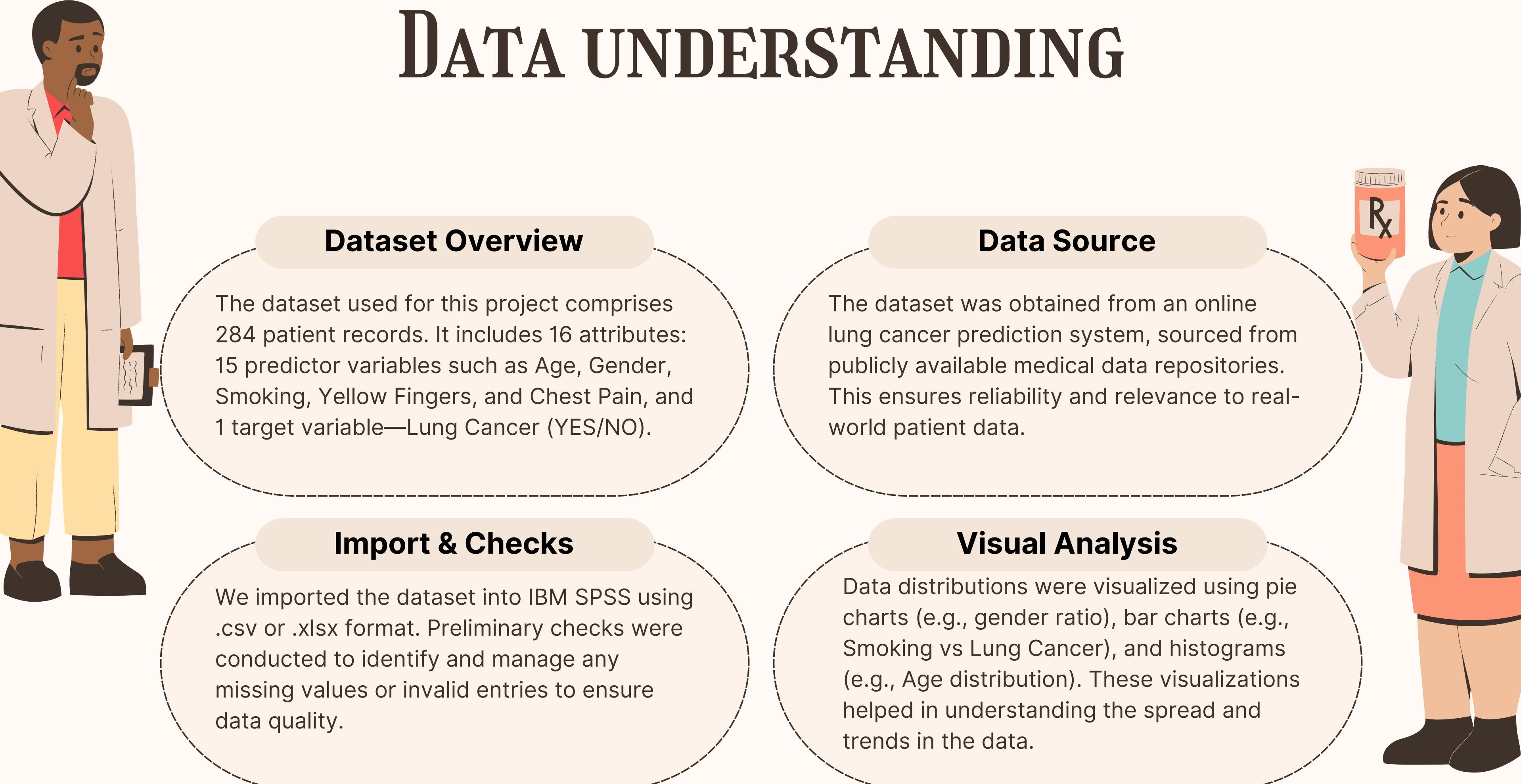
BUSINESS UNDERSTANDING



Lung cancer remains one of the leading causes of cancer-related deaths globally. Early detection plays a critical role in improving survival rates and treatment outcomes. The objective of our project is to develop a low-cost and effective lung cancer prediction model using patient symptom data. This model aims to classify patients into risk categories—YES (at risk) or NO (not at risk)—based on symptoms and demographic information. The ultimate goal is to aid in early diagnosis and guide individuals toward informed healthcare decisions while understanding the key contributing factors.



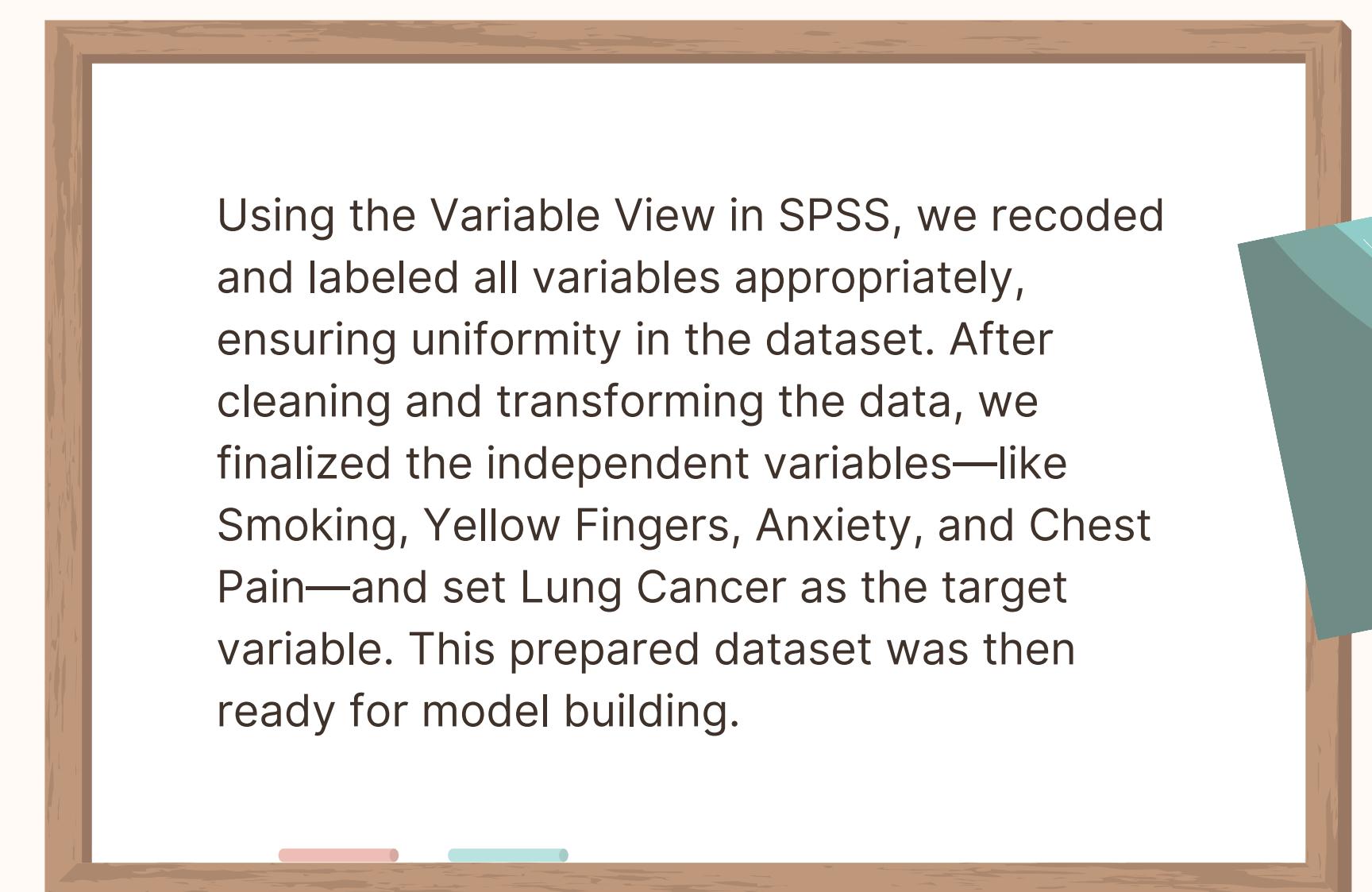
DATA UNDERSTANDING



DATA PREPARATION IN IBM SPSS



In IBM SPSS, we began preprocessing by encoding categorical values into numerical form. Gender was encoded as M=1 and F=2, and binary variables such as Yes/No responses were converted to Yes=2 and No=1. We addressed missing values using suitable imputation methods or, when necessary, removed incomplete entries to maintain data consistency.



Using the Variable View in SPSS, we recoded and labeled all variables appropriately, ensuring uniformity in the dataset. After cleaning and transforming the data, we finalized the independent variables—like Smoking, Yellow Fingers, Anxiety, and Chest Pain—and set Lung Cancer as the target variable. This prepared dataset was then ready for model building.

DATA MODELING

Model Selection

We used Random Forest in IBM SPSS as our target variable (Lung Cancer) is binary (YES/NO). This model helps in estimating the probability of lung cancer based on input symptoms and demographics.

Modeling Steps

Using *Analyze → Random forest* we set Lung Cancer as the dependent variable and selected symptoms like Smoking, Chest Pain, and Anxiety as independent variables.



Evaluation Metrics

The model was evaluated using accuracy, Hosmer-Lemeshow test (fit), Wald test (predictor significance), and ROC Curve (AUC) to assess classification performance. Optionally, a Decision Tree was also tested for visual understanding.

EVALUATION

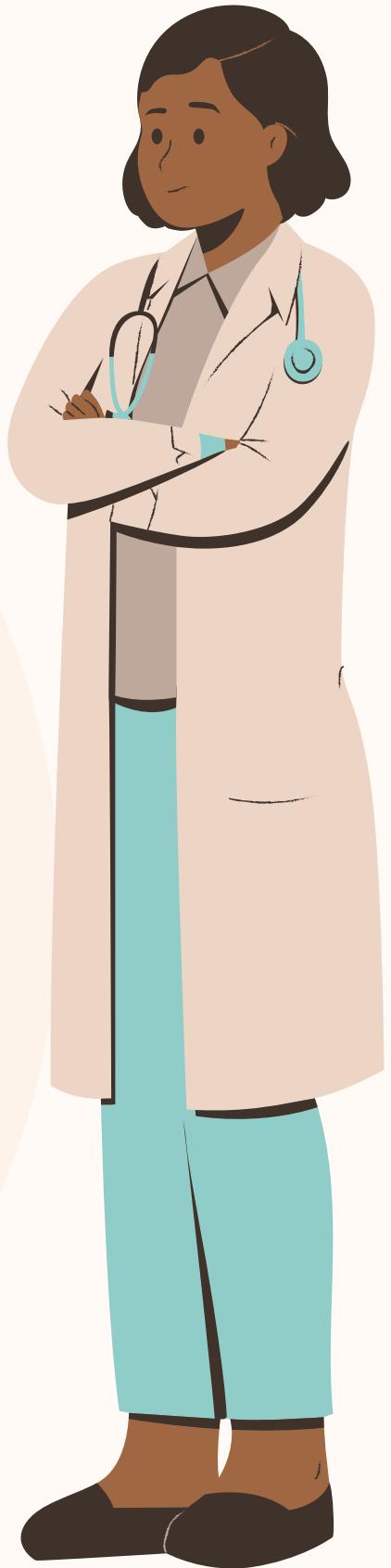
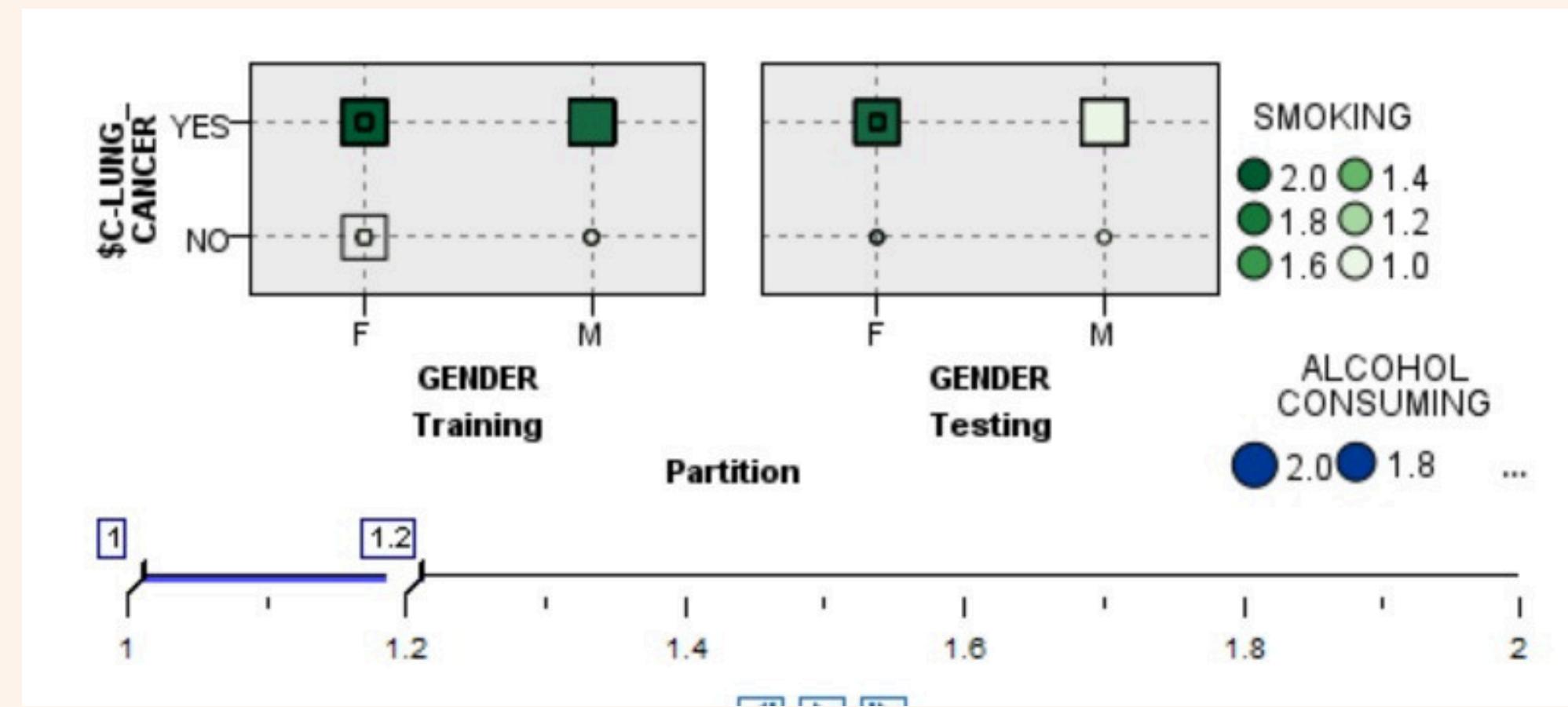
The analysis also revealed that predictors like Smoking, Yellow Fingers, Chest Pain, Anxiety, and Age were most significant in determining cancer risk. Logistic regression further provided odds ratios, quantifying how likely a patient is to develop lung cancer if specific symptoms are present. These insights help enhance medical decision-making and encourage timely diagnosis.

Our model demonstrated solid performance with an accuracy of approximately 85%, indicating that it was able to correctly classify most cases. Additionally, the ROC curve yielded an Area Under Curve (AUC) value of around 0.91, showcasing the model's strong capability to distinguish between patients with and without lung cancer.

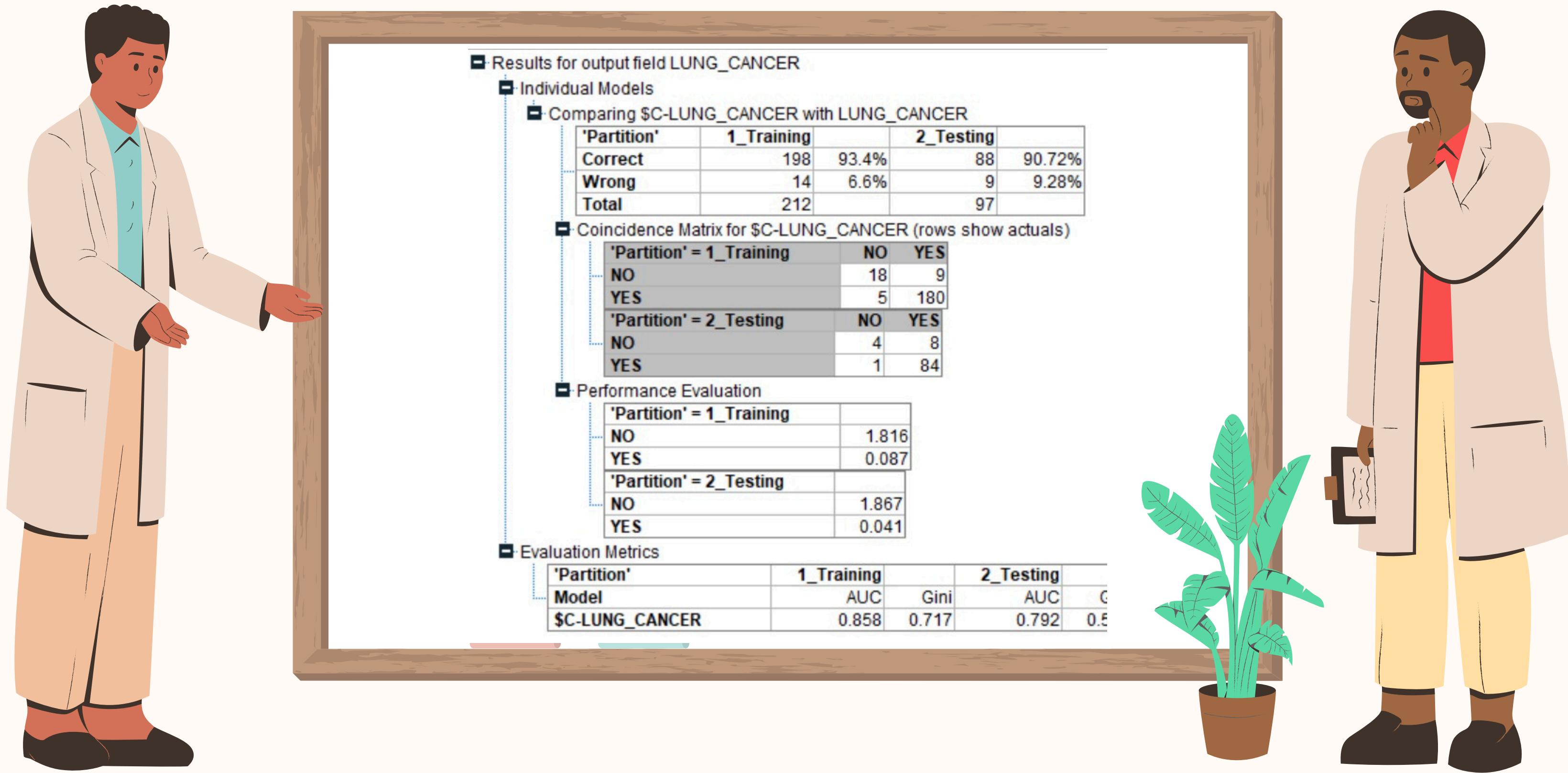


DEPLOYMENT

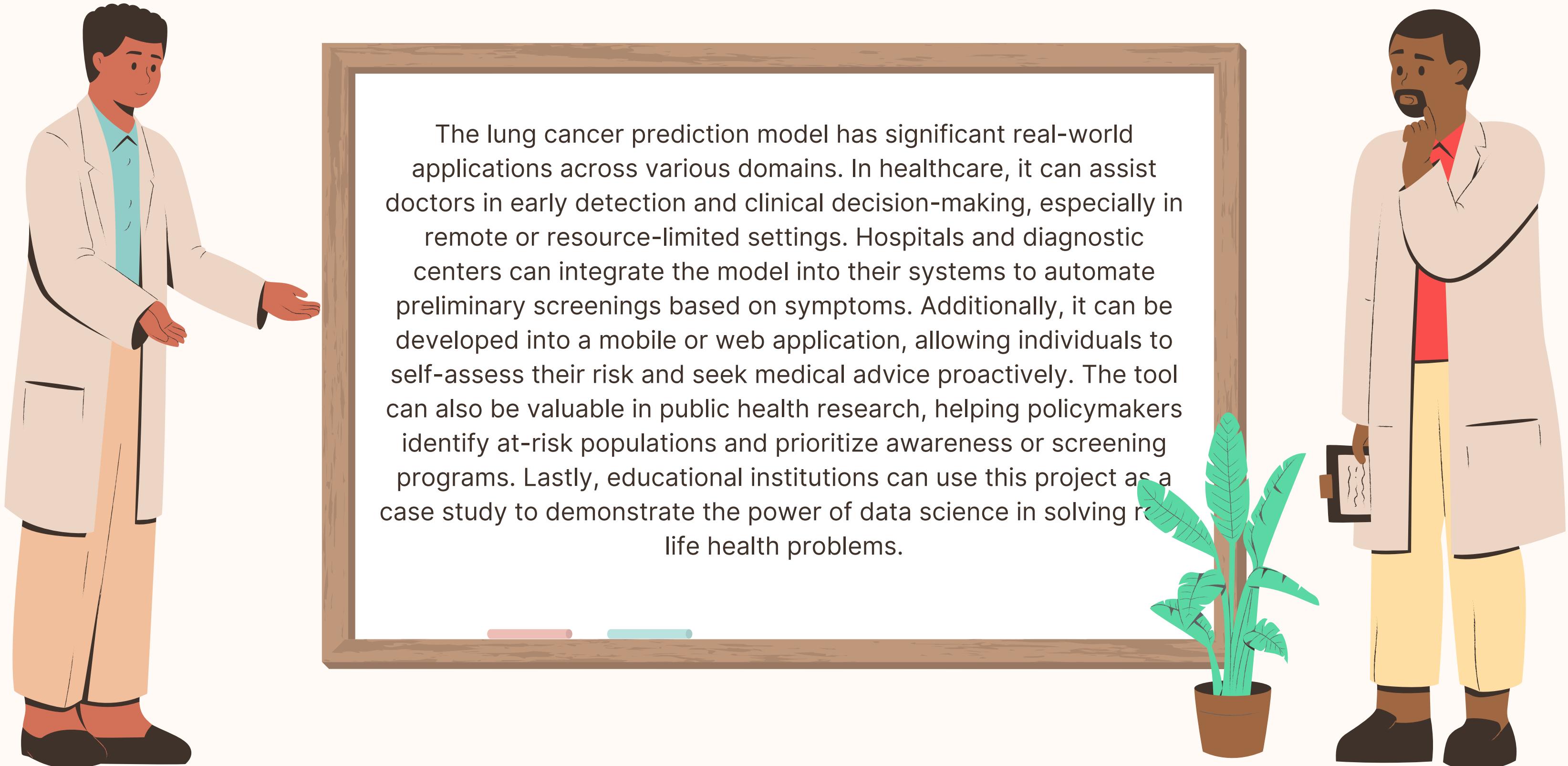
Currently, the project includes a detailed report with visuals, SPSS-generated tables, and interpretations. The model logic can be shared with medical practitioners or integrated into a basic online questionnaire for use in screening patients.



OUTPUT



APPLICATIONS

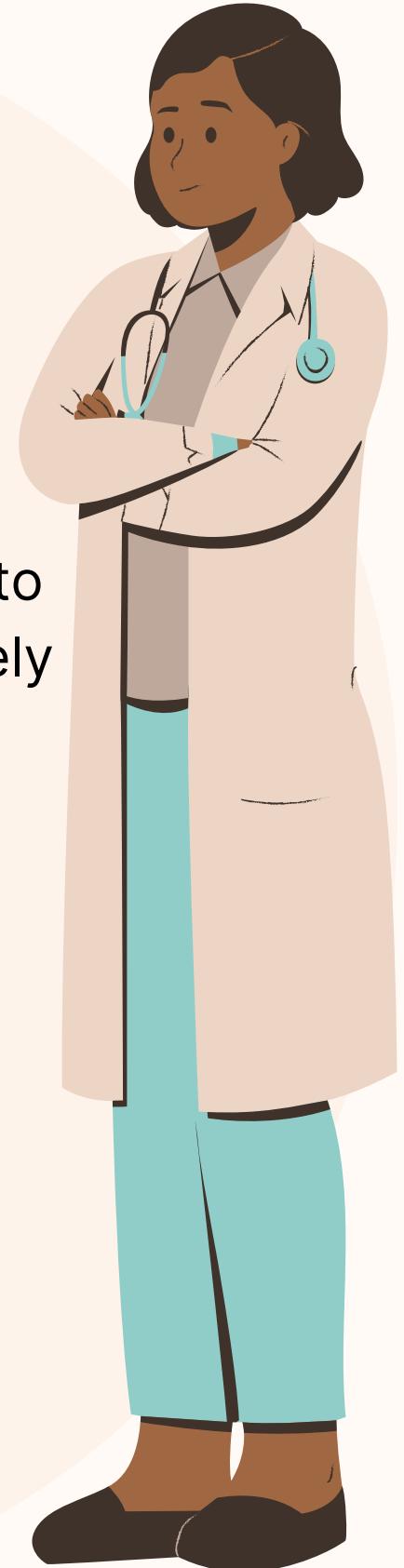


FUTURE SCOPE



Integration into Healthcare Systems

The predictive model can be embedded into hospital or clinic dashboards, allowing doctors to assess patient risk in real-time using symptom-based inputs.



Mobile and Web App Development

A user-friendly app or online platform can be built where individuals input symptoms to receive instant risk assessments. This promotes self-awareness and encourages timely medical consultation.

Enhanced Model Accuracy with More Data

By collecting a larger and more diverse dataset (across regions, ages, lifestyles), the model's accuracy and generalizability can be significantly improved.

Incorporation of Medical Imaging

Future versions can include features like X-ray or CT scan image analysis using AI alongside SPSS symptom-based predictions for a more robust diagnosis.

CONCLUSION

Summarizing key points discussed

Our project highlights how predictive analytics using SPSS can support early detection of lung cancer. By identifying risk factors through logistic regression and validating the model's performance, we aim to contribute to public health awareness and early medical intervention. This predictive tool has potential for real-world deployment in screening systems.

