

12.9.23

MTA School of Computer Science  
Project Number 231006

## ADVISORS

Prof. Adi Shraibman  
Dr. Dorit Shweiki  
Dr. Yonatan Bilu

## STUDENSTS

Dan Naftaly  
Tamir Matok

# Medical Data Analysis

**lung Cancer Research**



Table of Contents



		Page
I	Research Background & Motivation	3
II	Exploratory Data Analysis	6
III	Model & Techniques	8
IV	Model evaluation	10
V	Conclusion	14

## Research Question

While smoking is the leading cause of lung cancer, many non-smokers also develop the disease. Other possible causes include exposure to air pollution, radon, asbestos, and genetic factors.

We want to rank and investigate the different risk factors and understand how they contribute to the development of lung cancer.



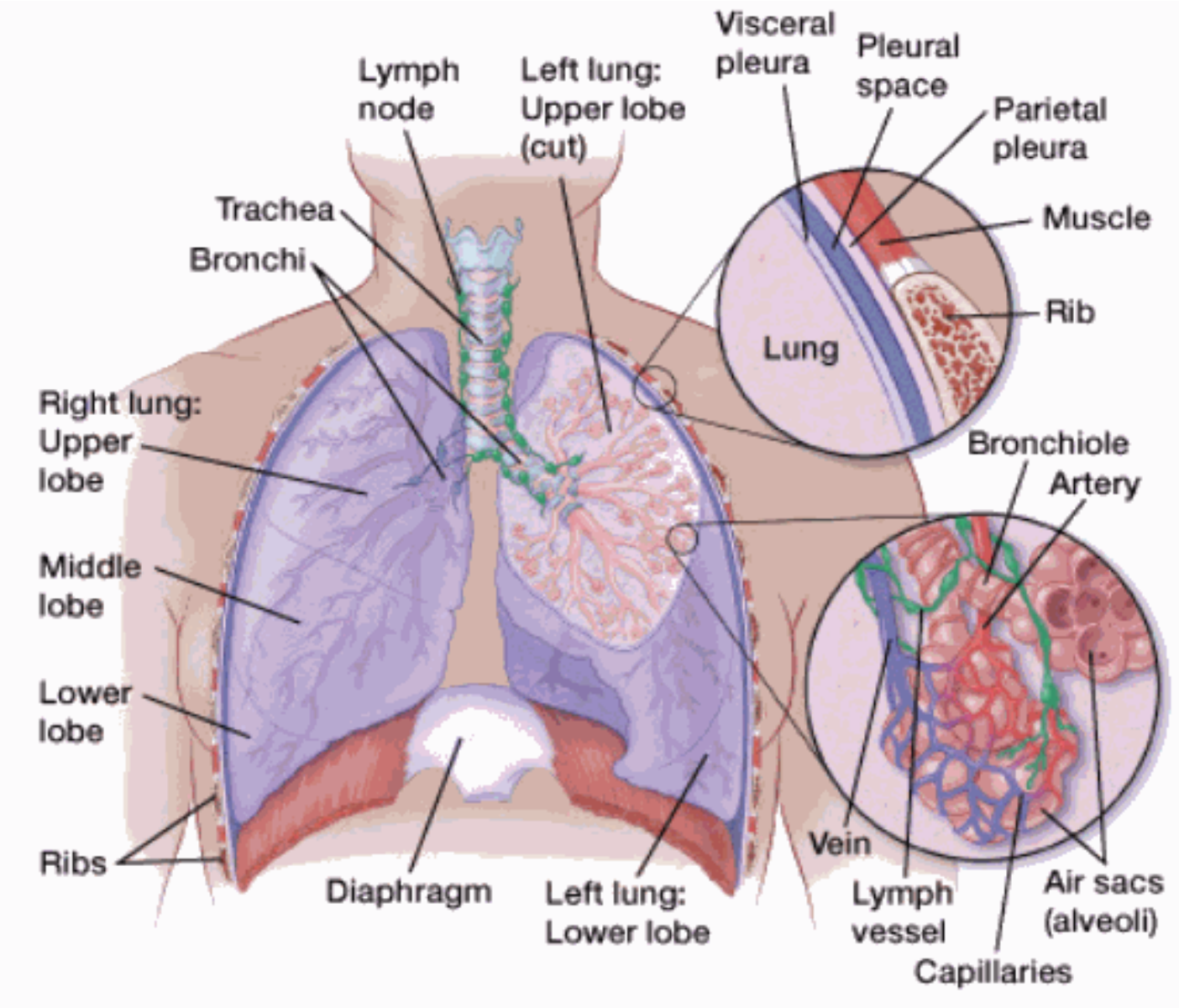
## Motivation

- Personal connection as smokers adds a layer of personal significance to this research, driving us to unravel the complexities surrounding lung cancer risks.
- Need to understand the diverse factors contributing to the development of lung cancer.
- Uncovering insights into potential contributors such as air pollution, radon, asbestos, and genetic predisposition, with the ultimate goal of enhancing prevention and treatment strategies for this devastating disease.





# Selected types of lung cancer



C34	Malignant neoplasm of bronchus and lung
C34.0	Main bronchus Carina Hilus (of lung)
C34.1	Upper lobe, bronchus or lung
C34.2	Middle lobe, bronchus or lung
C34.3	Lower lobe, bronchus or lung
C34.8	Overlapping lesion of bronchus and lung [See note 5 at the beginning of this chapter]
C34.9	Bronchus or lung, unspecified

# Data Fields

## Behavioral

- Smoking status
- Number of cigarettes previously smoked daily
- Chest pain during physical activity
- Doctor diagnosed asbestosis
- Doctor diagnosed COPD
- Cough on most days

## Physiological

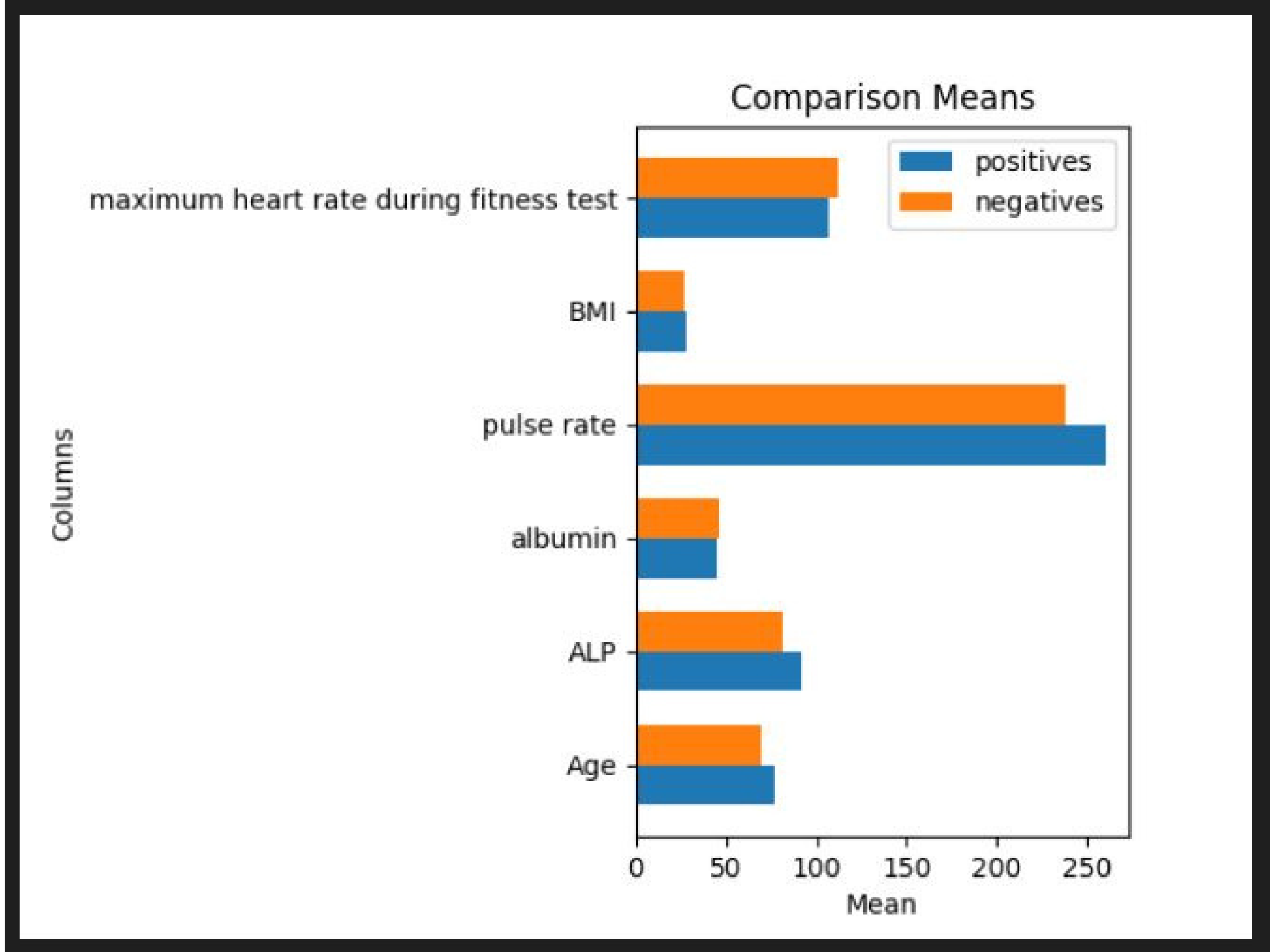
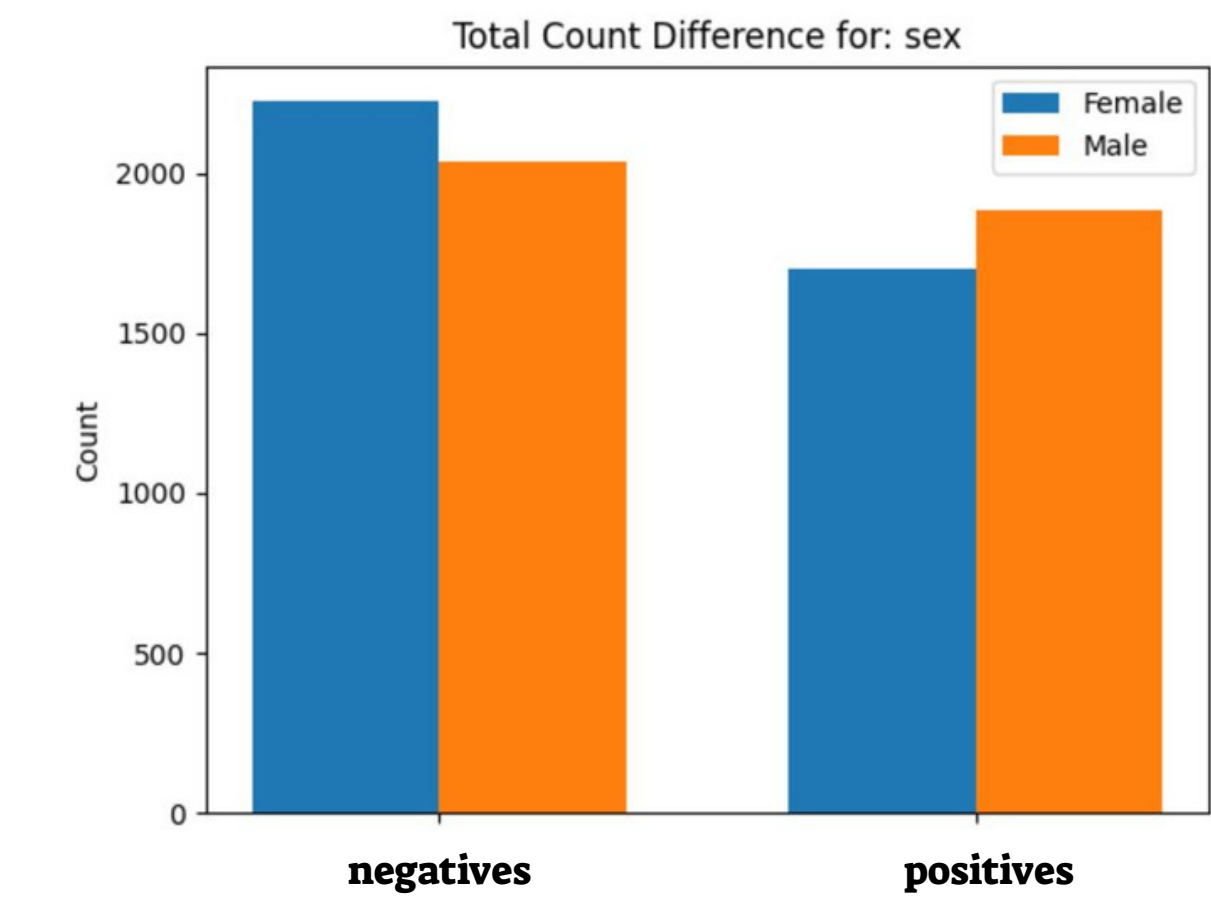
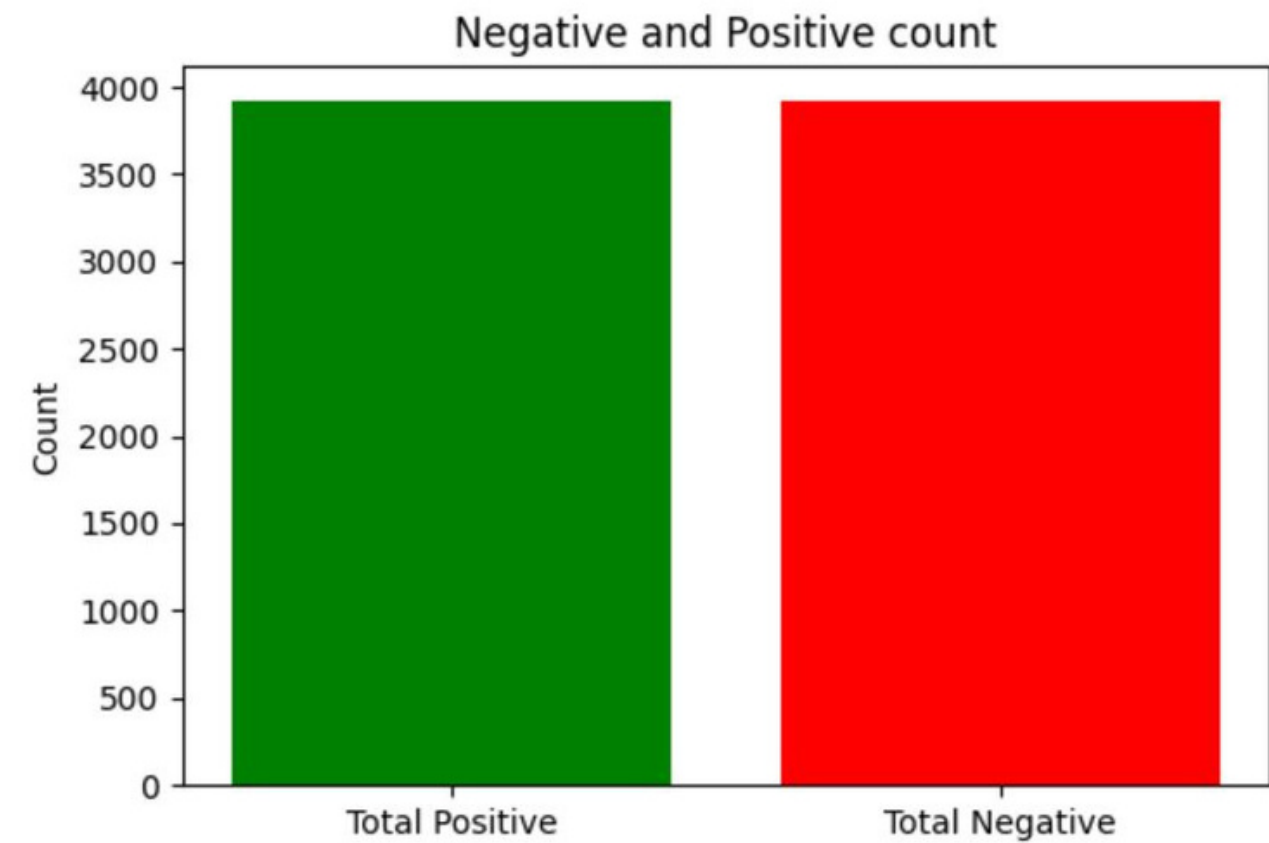
- |                    |  |
|--------------------|--|
| Sex                | Lymphocyte count                       |
| Year of birth      | ALP                                    |
| BMI                | White blood cell count                 |
| Albumin            | Forced Vital Capacity                  |
| Platelet count     | FEV1                                   |
| Pulse rate         | Maximum heart rate during fitness test |
| C-reactive protein |  |

## Demographic

- Ethnic background

## Environmental

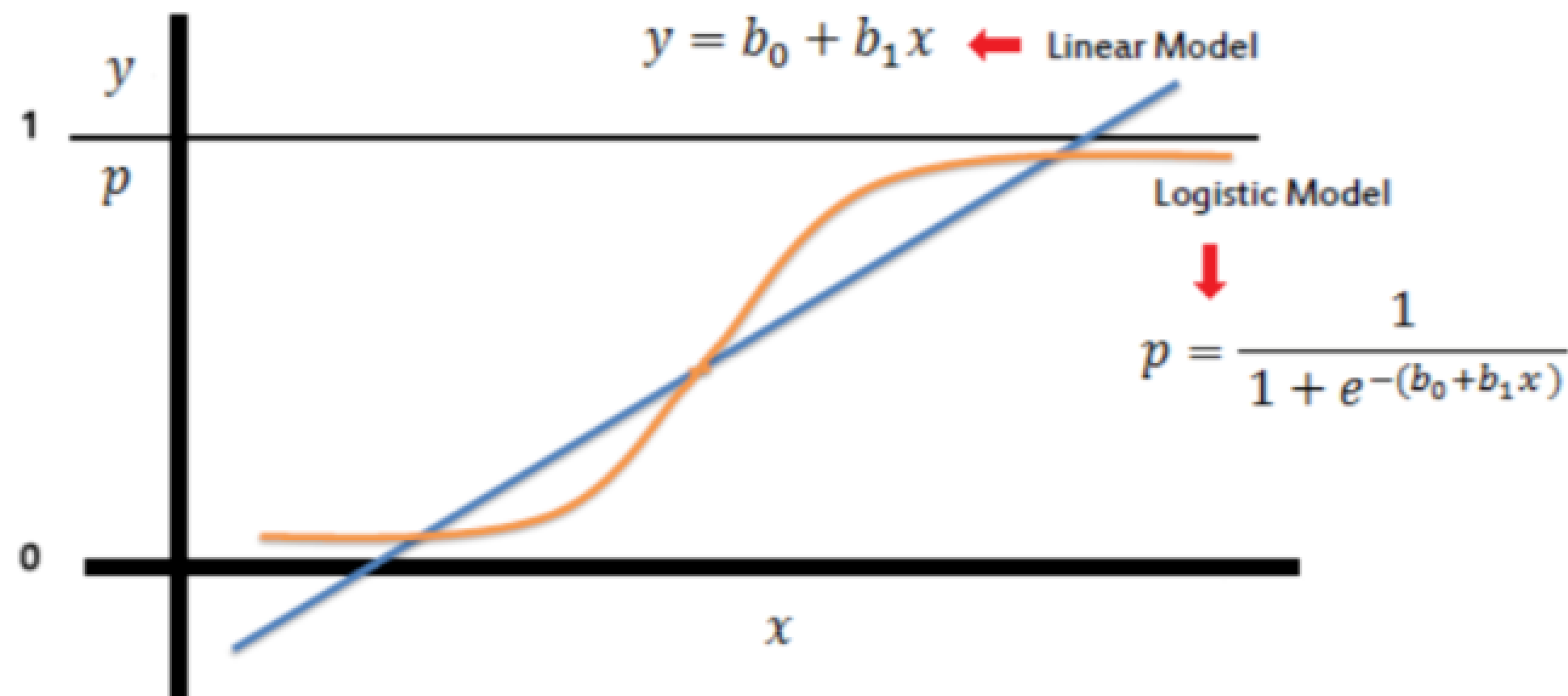
- Close to major road
- Nitrogen dioxide air pollution
- Particulate matter air pollution 2.5-10um



# Model selection

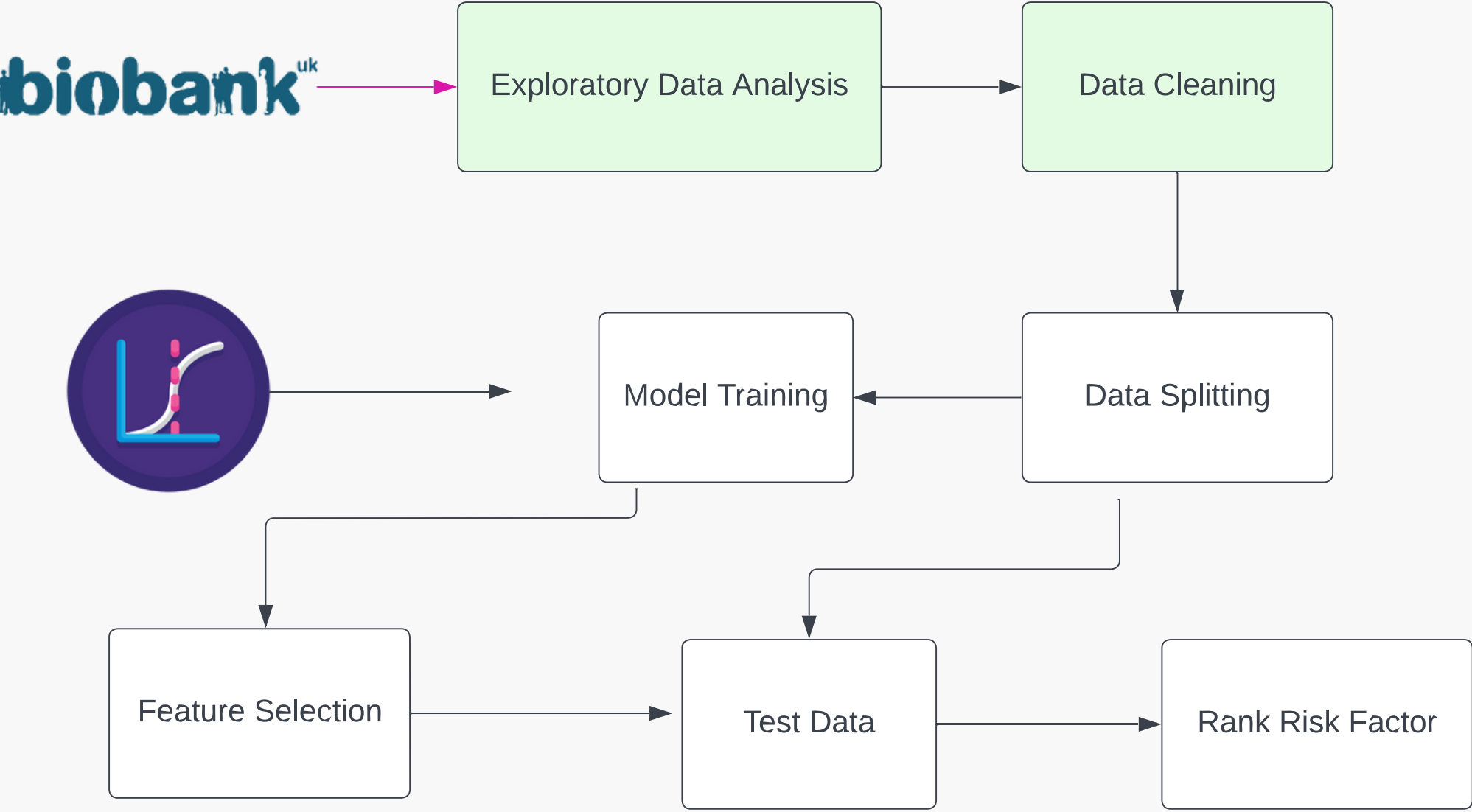
## Logistic Regression

Logistic regression is a supervised learning algorithm that makes use of logistic functions to predict the probability of a binary outcome.

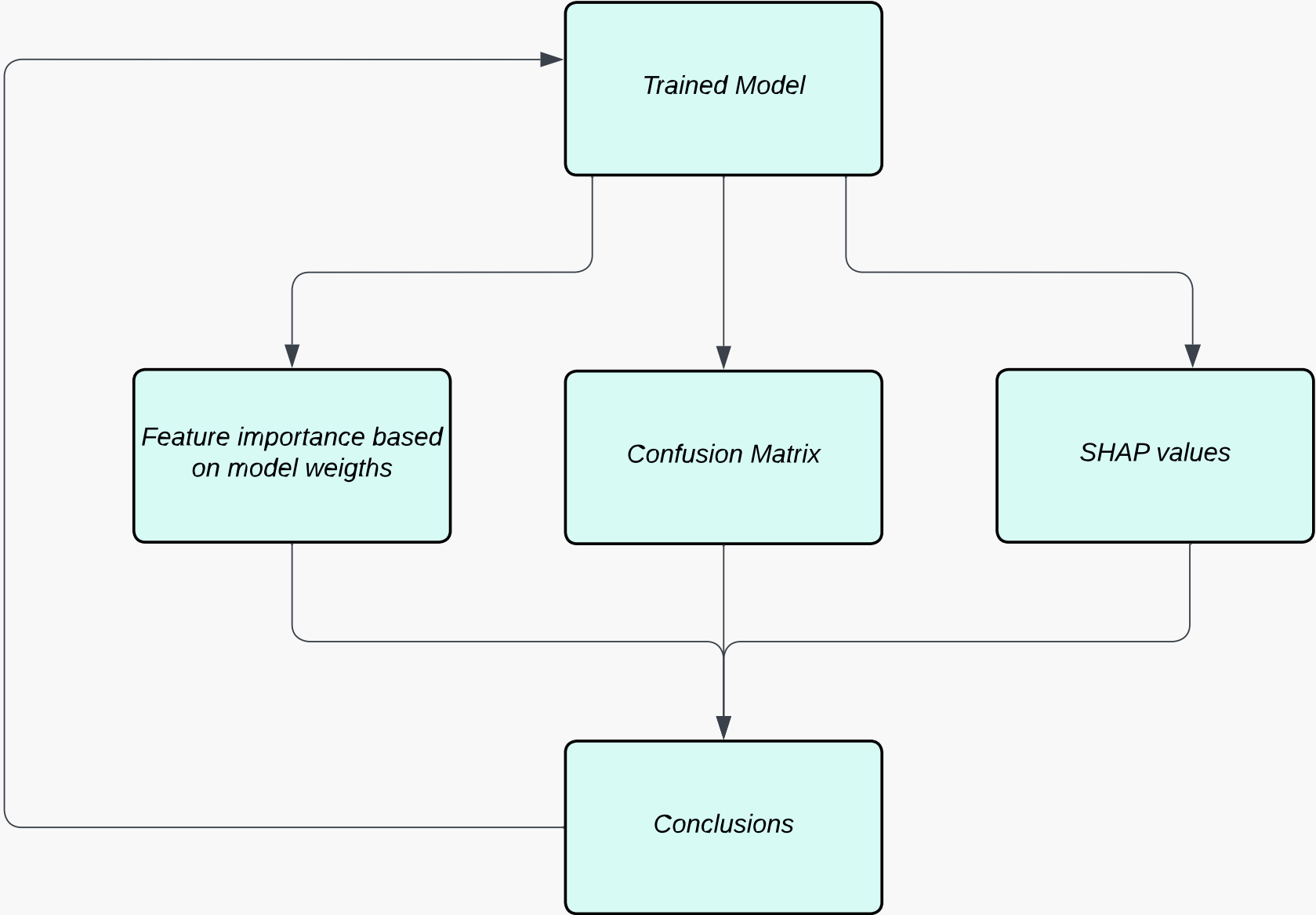


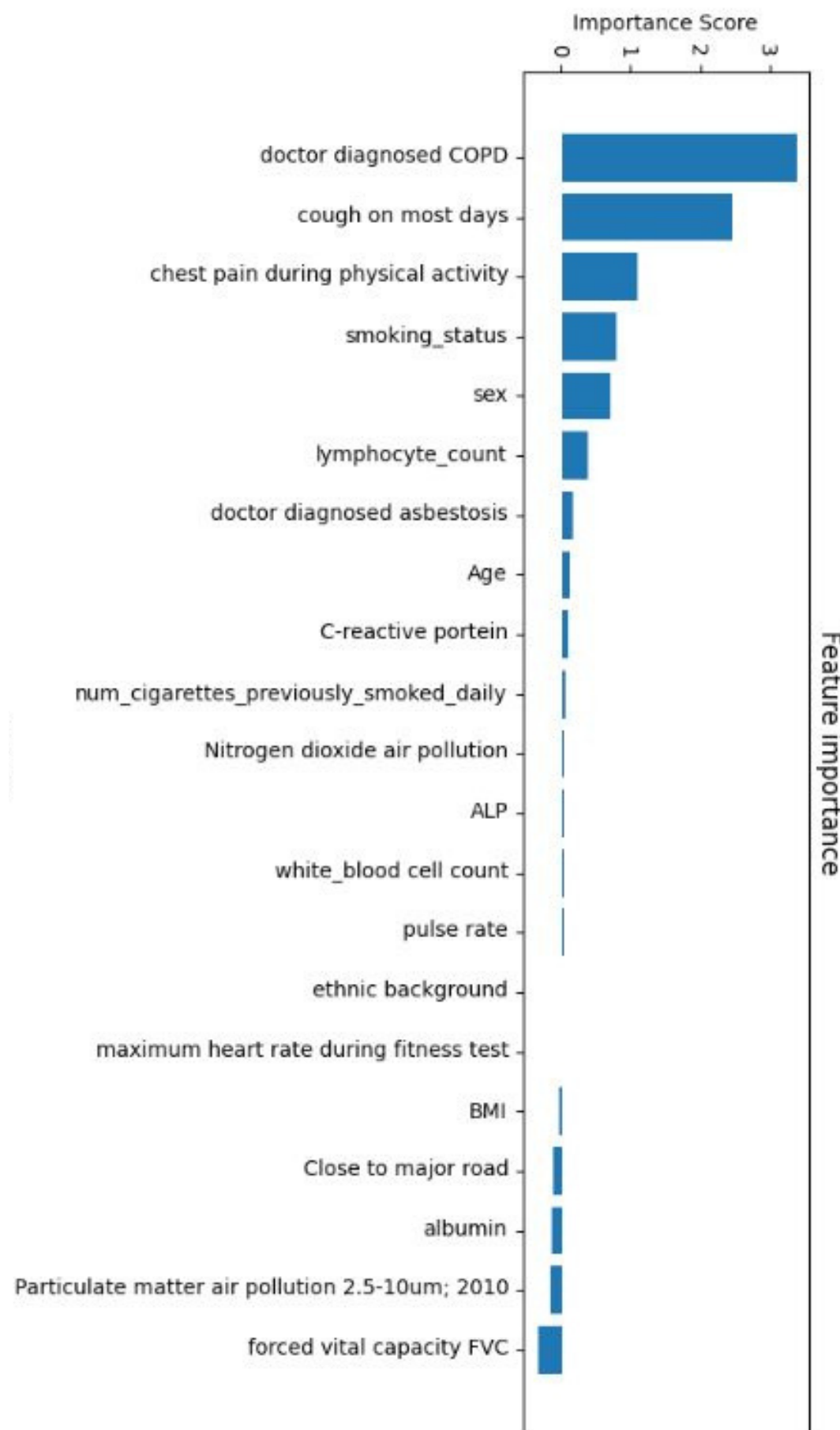


# High Level Architecture and Design

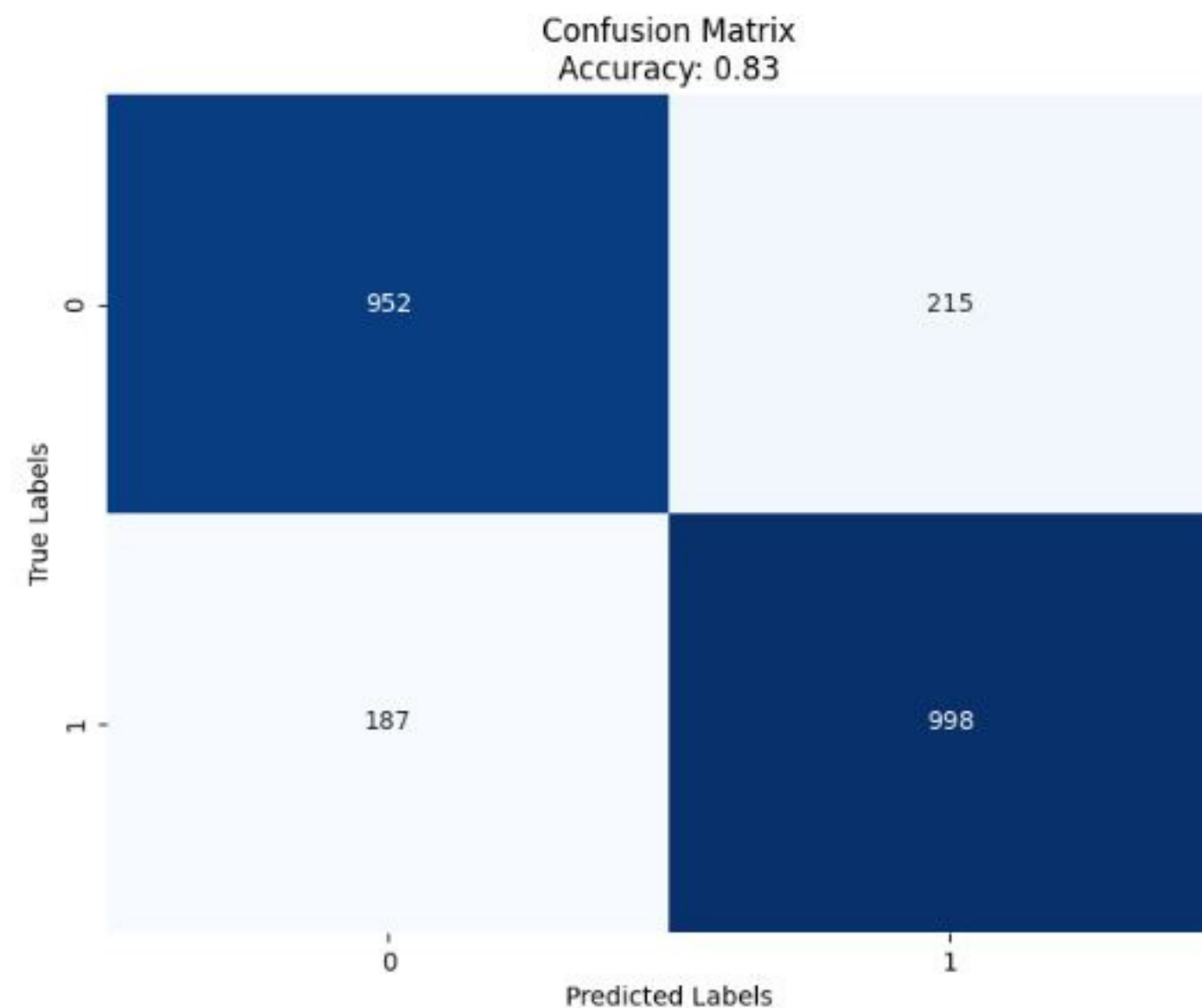


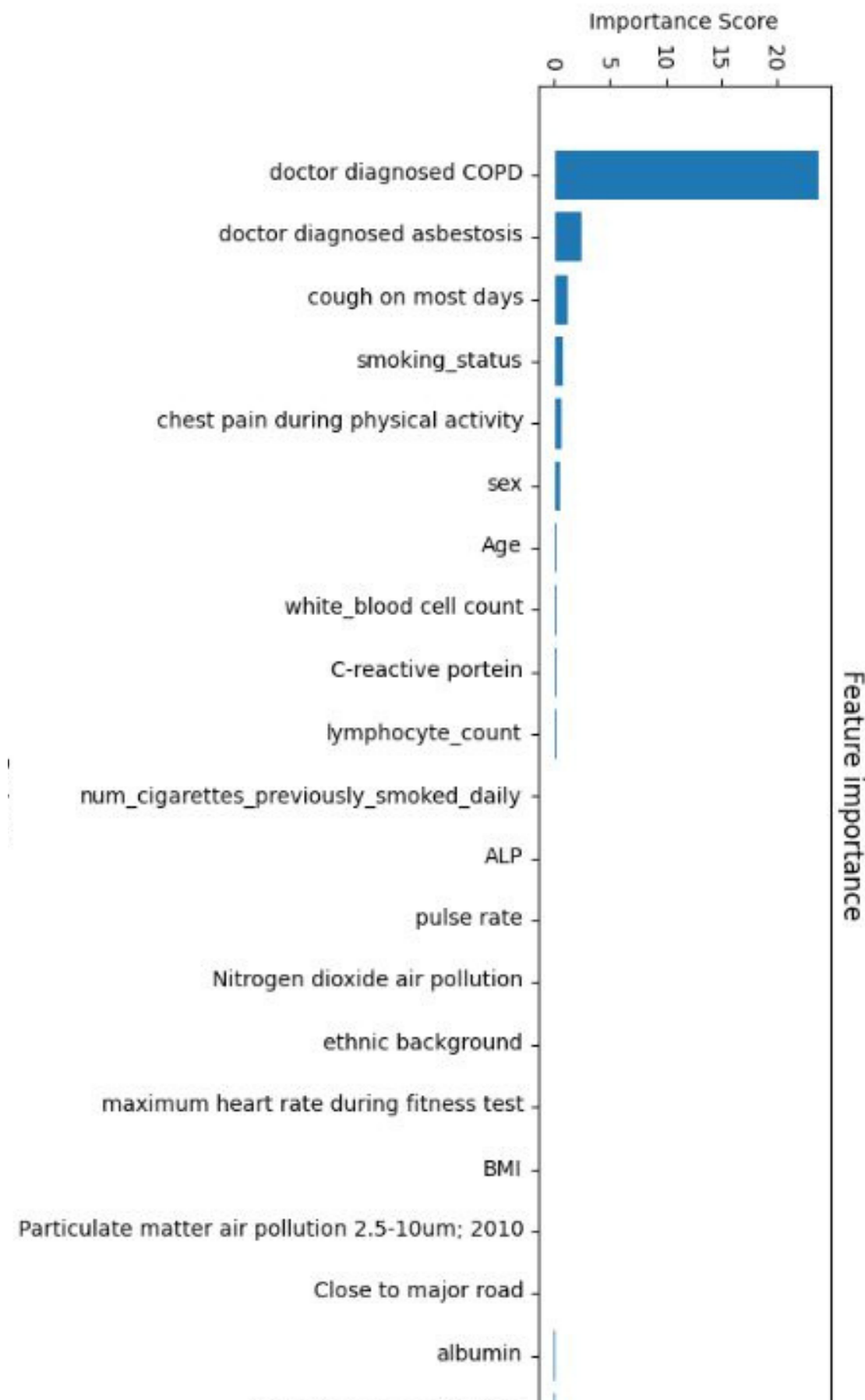
# Model evaluation process



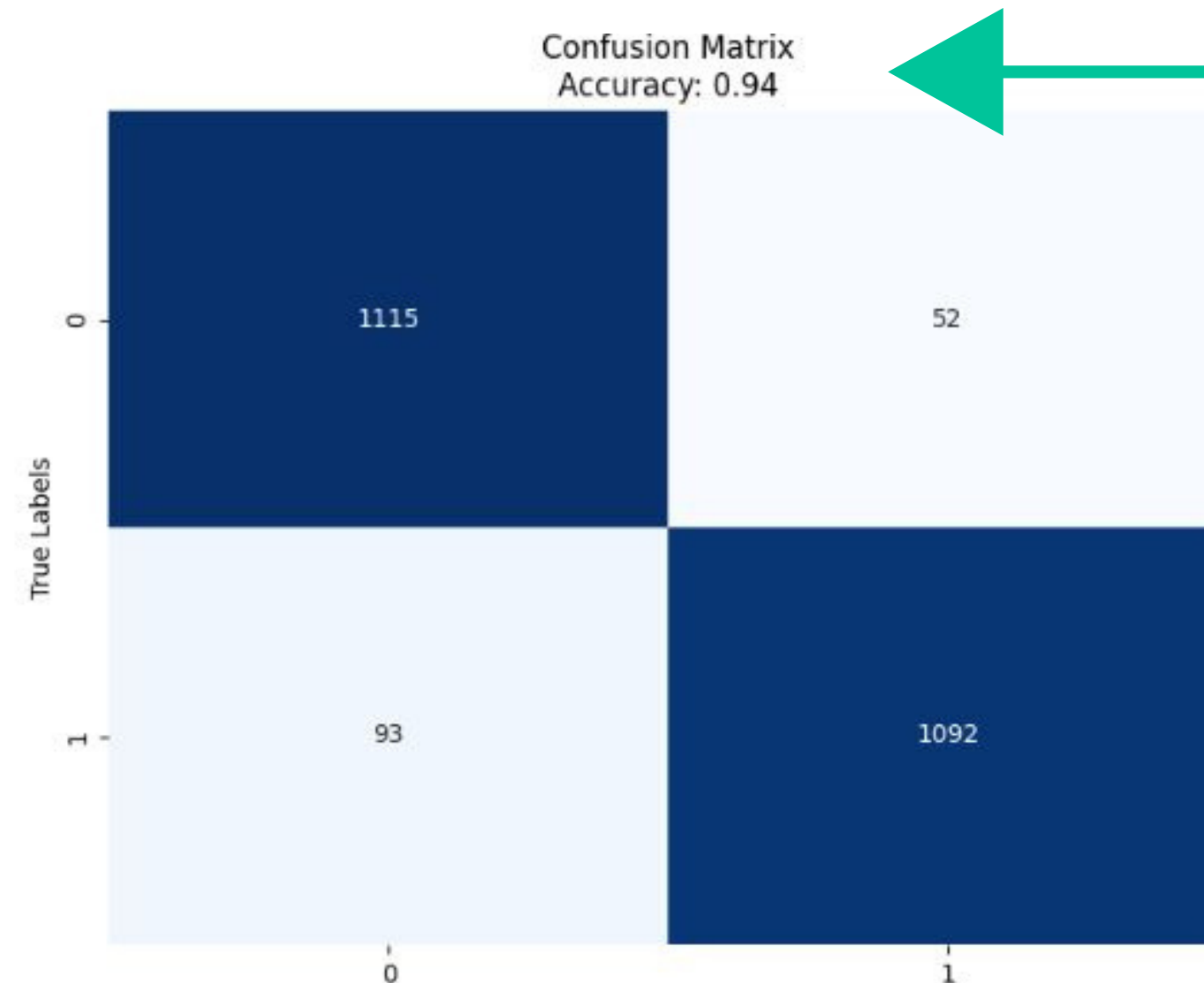


- **At first glance - model results are satisfying based on the confusion matrix.**
- **Several columns seem to have a significant impact on the model results, which raises concerns about the authenticity of our results.**

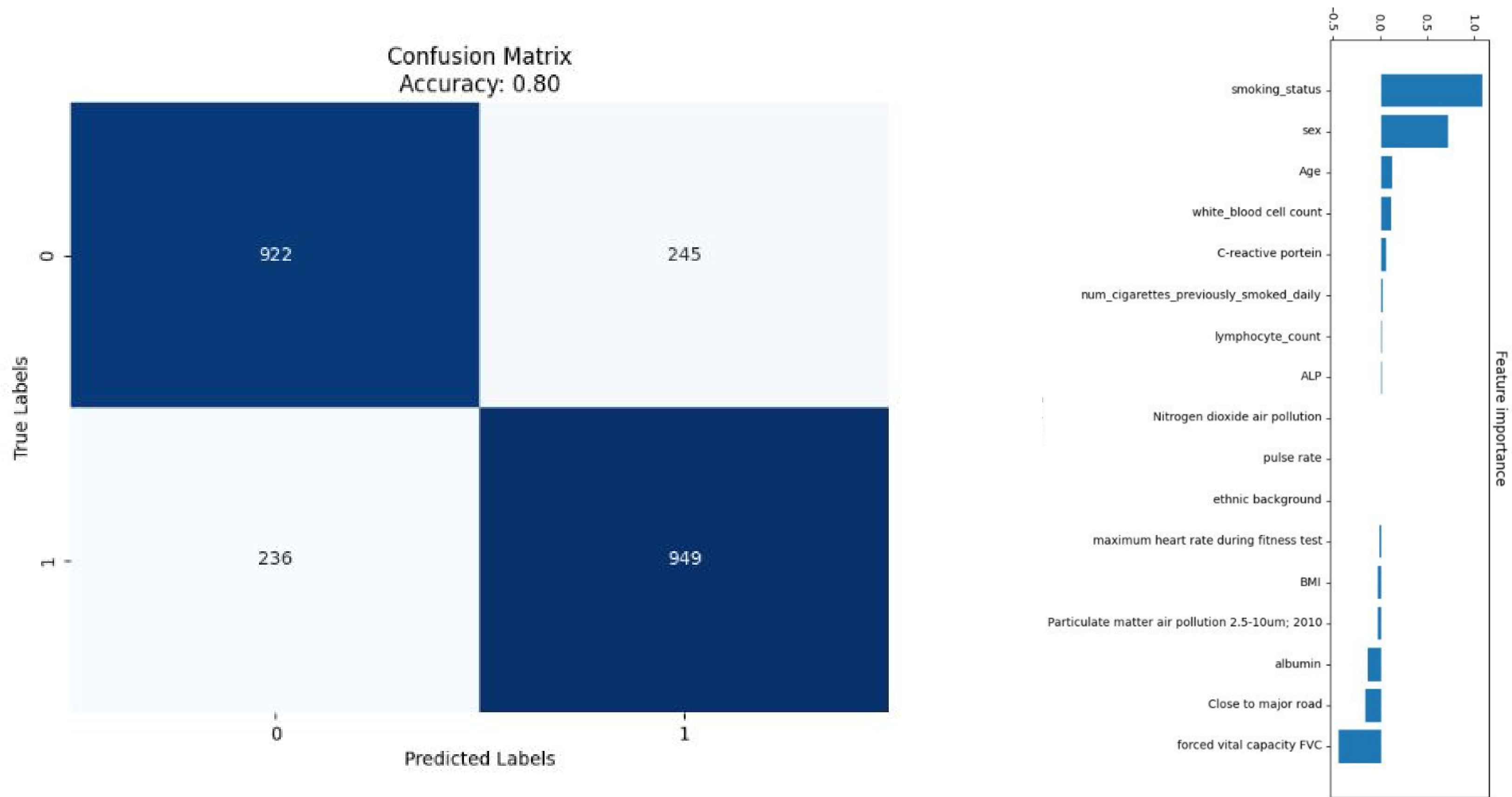




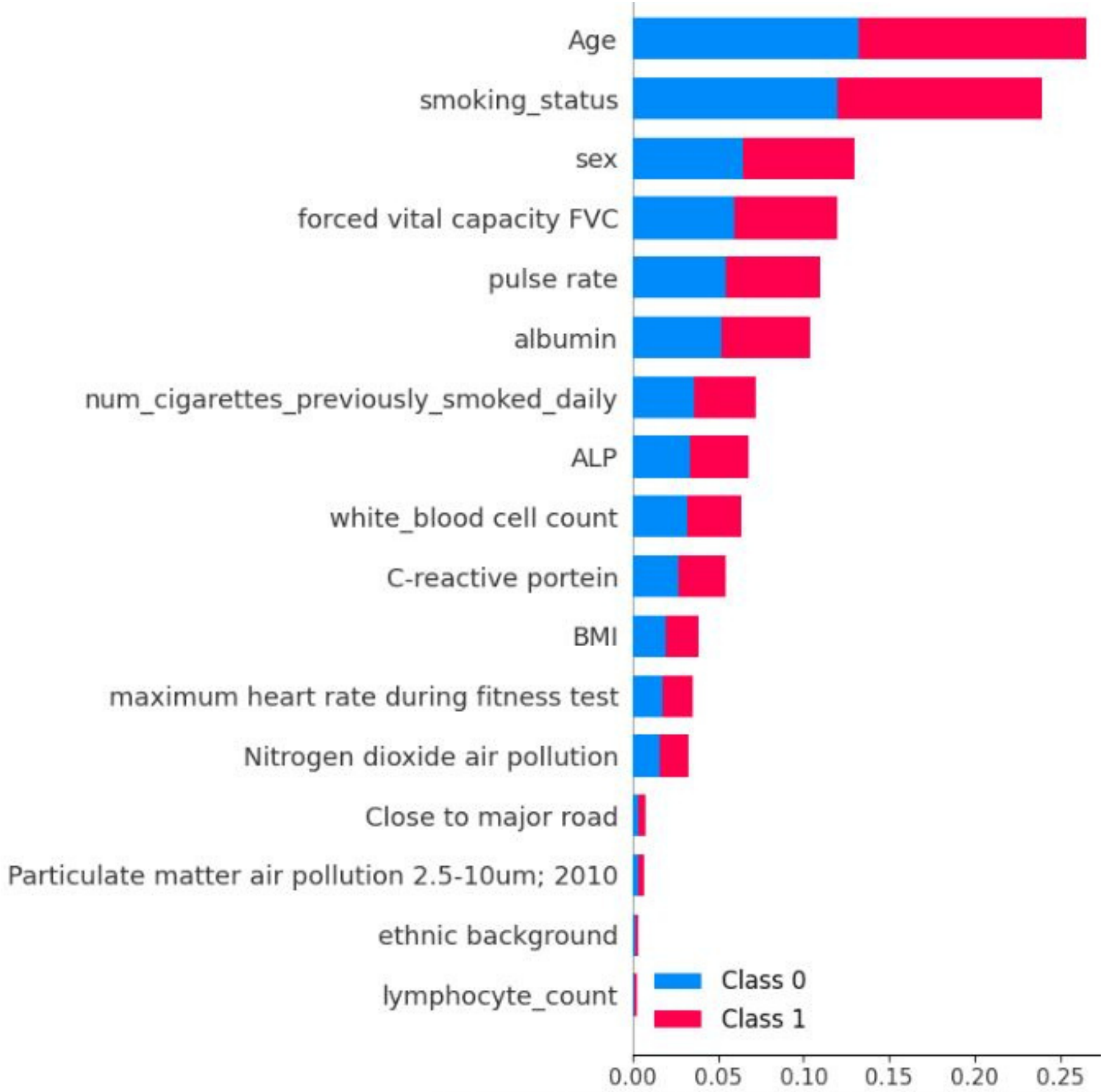
- By incorporating L1 regularization, we can identify and select the most important features related to lung cancer risk.
- Among smokers, COPD is an important risk factor for developing lung cancer, and predates lung cancer in up to 70–80% of cases (National Institutes of Health)



# Final Model Results



# SHAP values





# Summary

- The model has provided valuable insights into the factors influencing this deadly disease.
- With age and smoking status emerging as the foremost contributors to lung cancer risk, the findings reaffirm the importance of addressing these well-known risk factors.
- Exploration of additional variables, such as ethnic background and air pollution, revealed their comparatively lower significance in the context of lung cancer.

This research not only enhances our understanding of the disease but also underscores the need for continued efforts in smoking cessation and age-appropriate screenings to combat lung cancer effectively.