

Scientific Visualization Final Project

SNF-Driven Patient Stratification and Vital Status Prediction

Nada Bou Kanaan

Datasets

Clinical data Adenosquamous carcinoma patients: clinical data of 138 patients with Adenosquamous carcinoma with 24 features: age, sex, race, participant_country, tumor_site....additioal_pathologic_findings, bmi, alcohol_consumption, tobacco_smoking_history, follow_up_days, 'vital_status'..

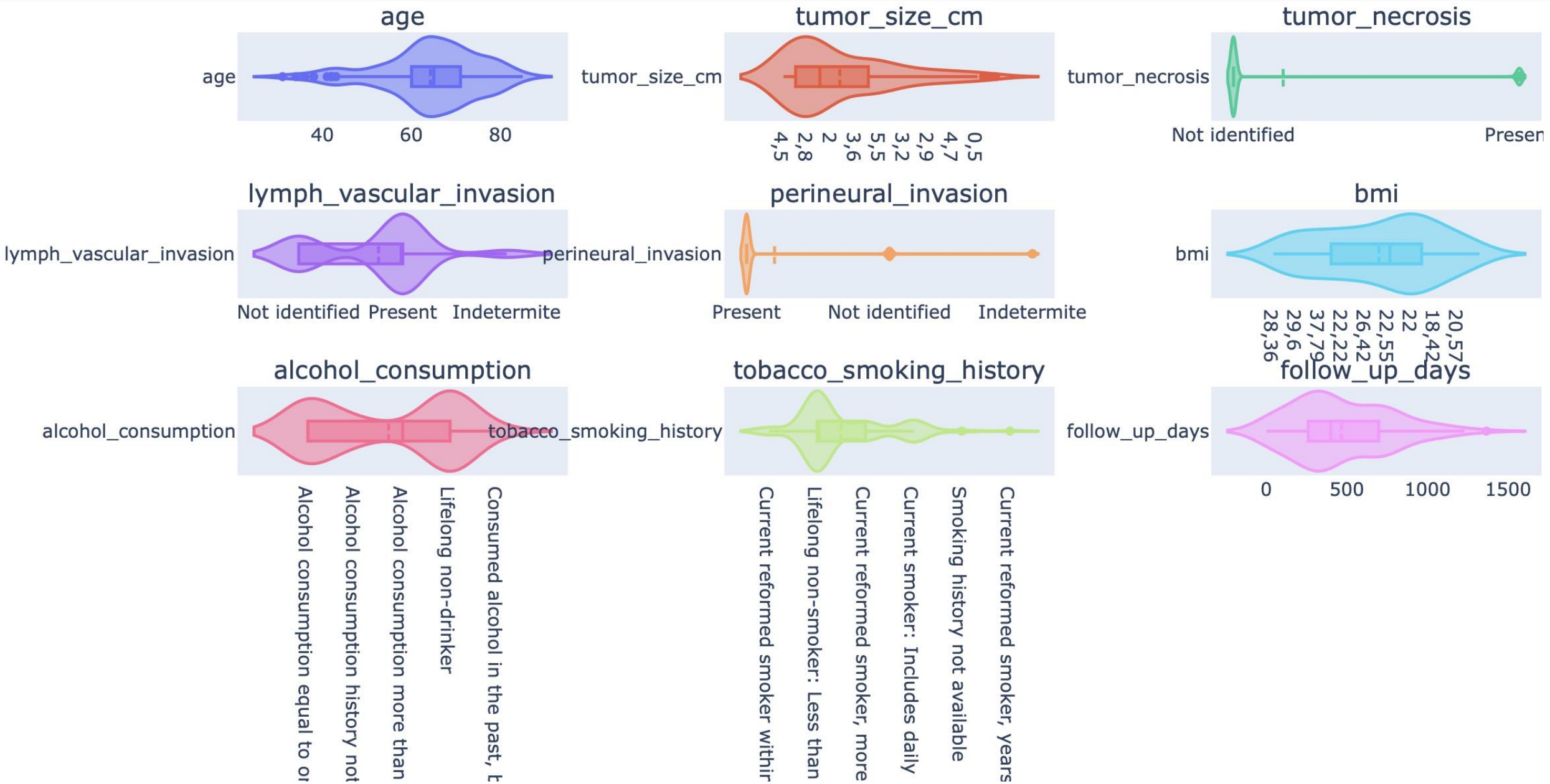
SCNA_gene: gene level measurement of somatic copy number alterations for the same 138 patients

mRNA: mRNA expression data for the same 138 patients

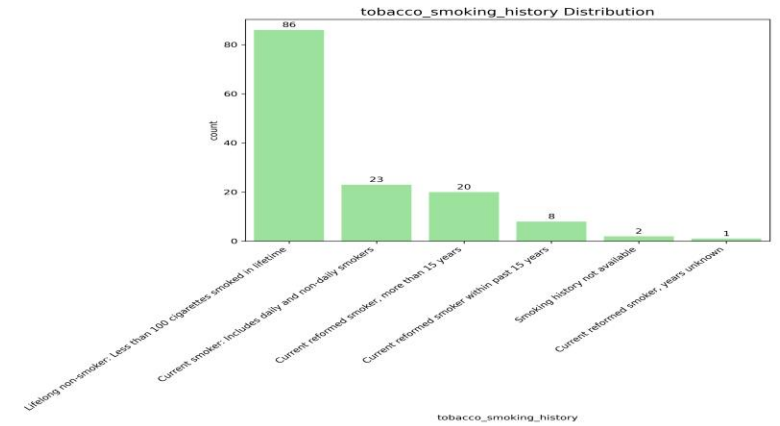
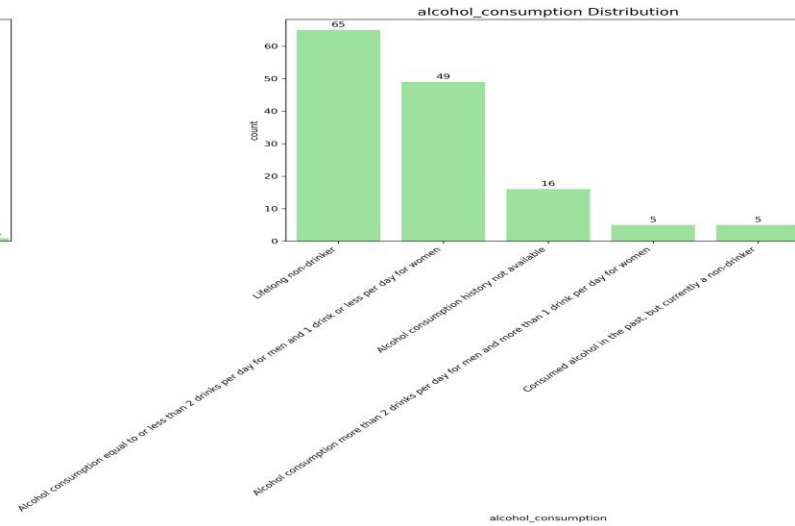
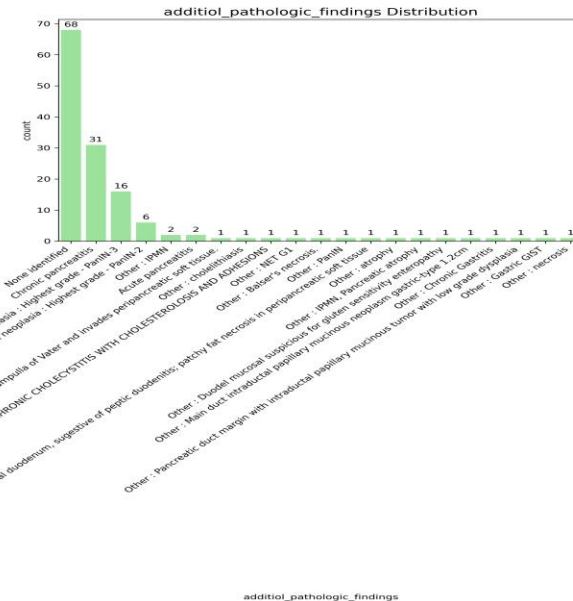
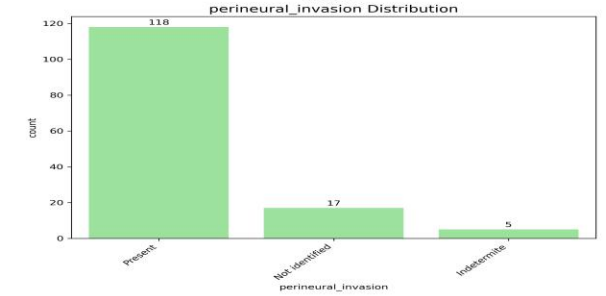
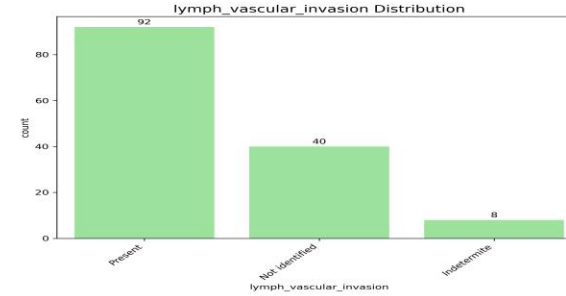
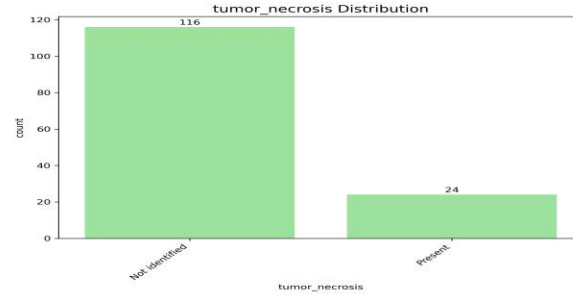
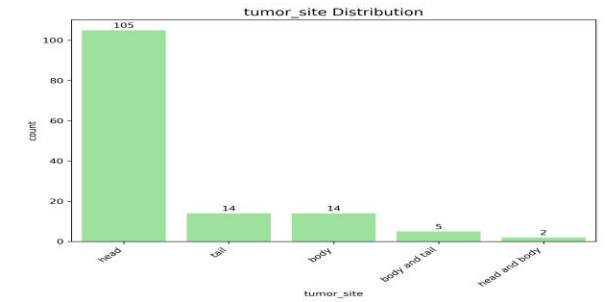
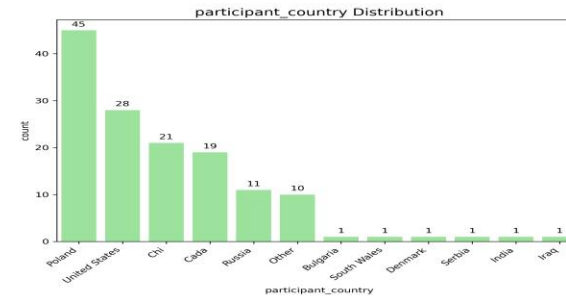
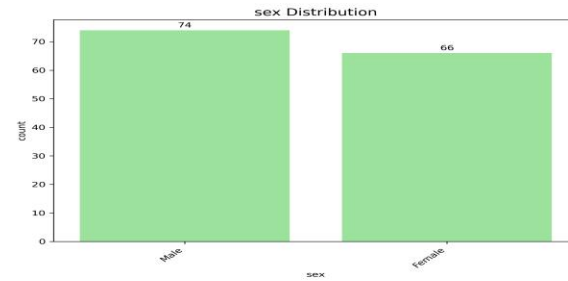
First step: prediction without the genomics data

- This code analyzes the clinical data of 138 patients with Adenosquamous carcinoma. I trained two classifiers—Random Forest and K-Nearest Neighbors (KNN)—to predict the patients' vital status based on their clinical data. Finally, I used the models to predict the vital status for the four patients with missing information. I visualized the data composition and the feature importance to better understand the distribution of the clinical variables and which features most influenced the predictions.

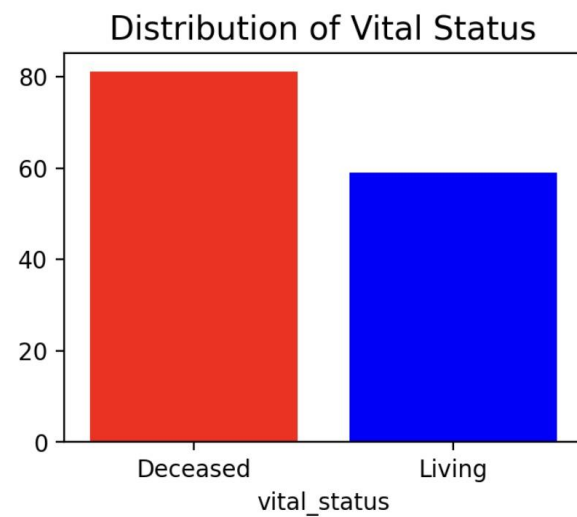
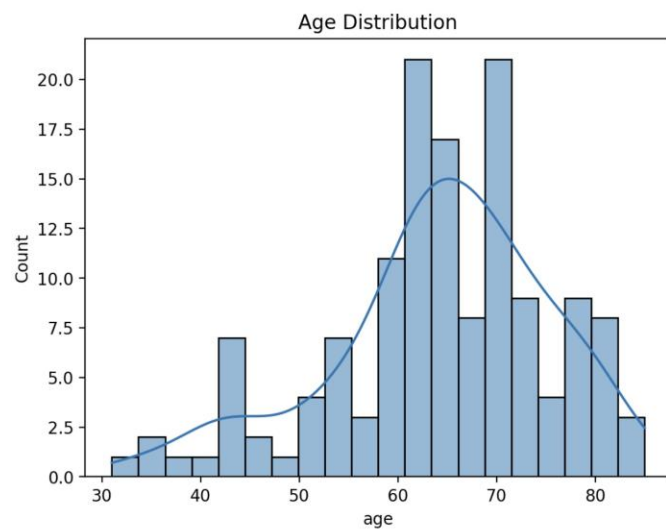
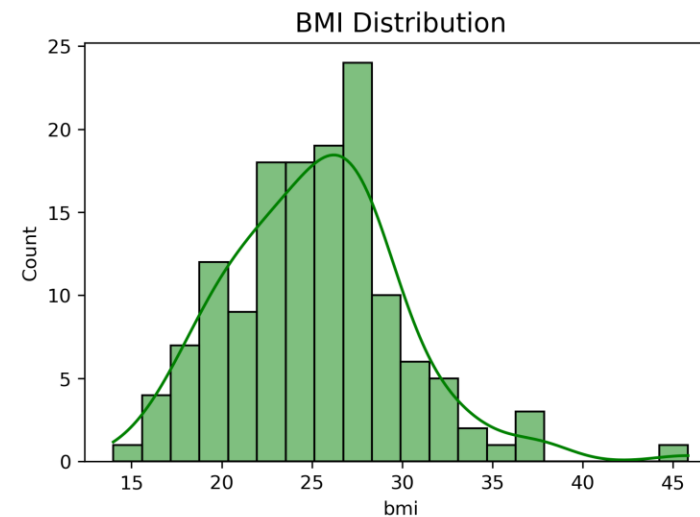
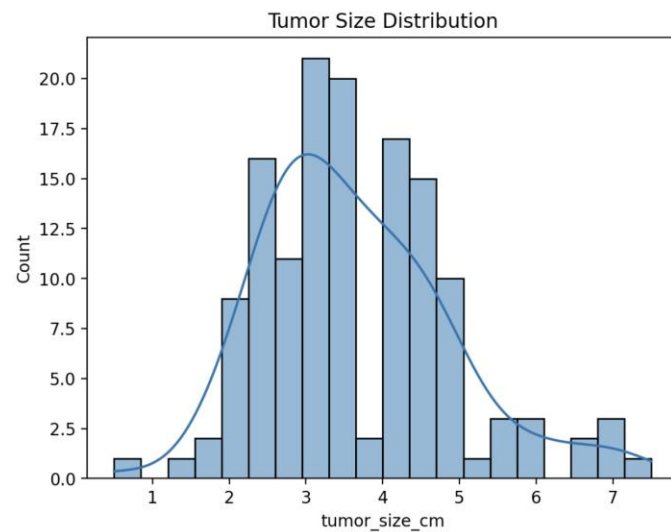
Data analysing



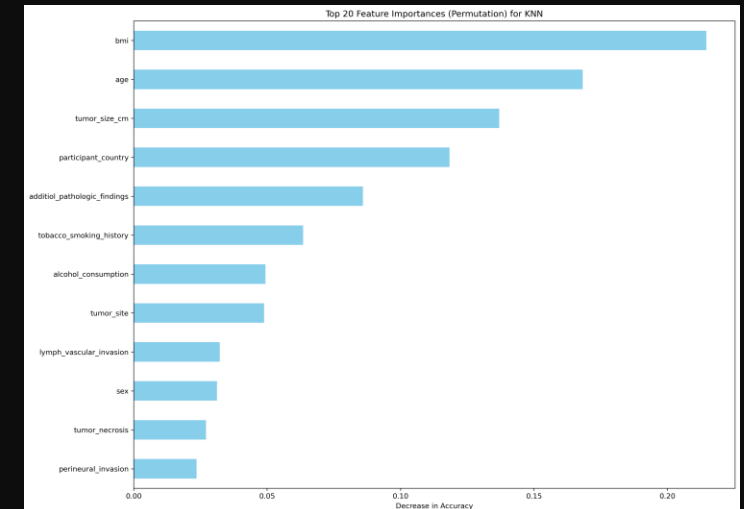
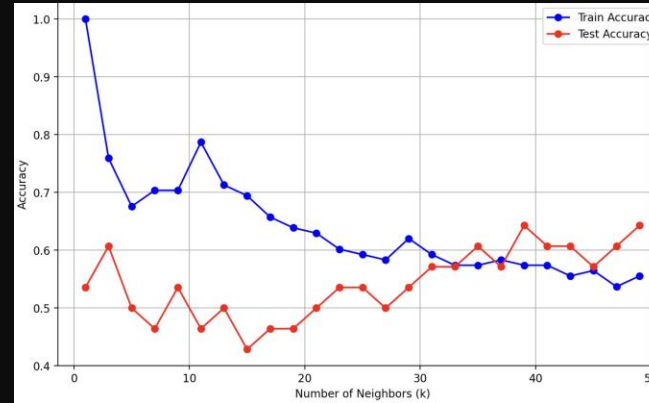
Data analysing



The world map displays the number of countries with which each country has trade relations. The color scale ranges from 0 (dark purple) to 35 (green). The United States and Canada are colored dark blue, indicating approximately 20-25 trade partners. Most European countries are colored green, indicating 30 or more trade partners. China and Japan are colored dark blue, indicating approximately 20-25 trade partners. India is colored dark purple, indicating 0 trade partners. Many countries in Africa and South America are also colored dark purple, indicating 0 trade partners. The map shows a high concentration of trade relations in North America, Europe, and East Asia.



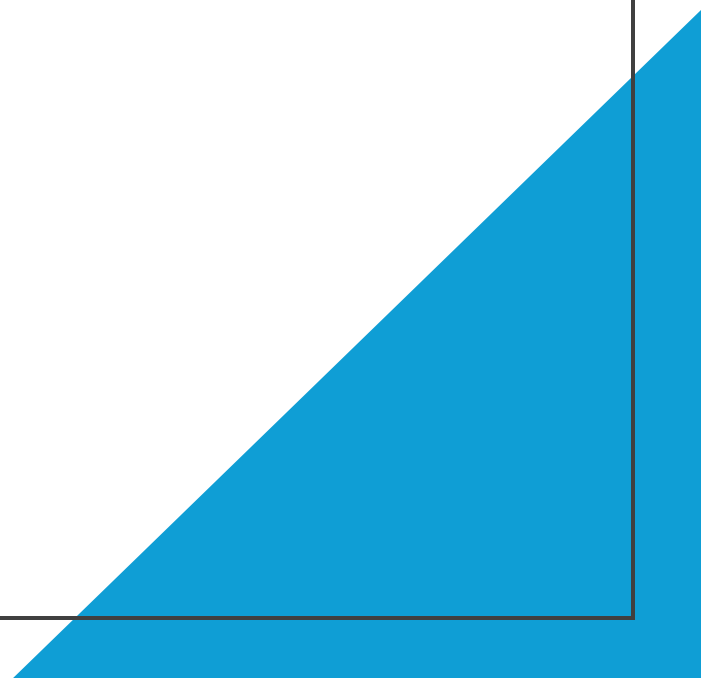
KNN feature
importance: bmi, age
optimal number of
neighbors : best k 39



Predictions

- | case_id | vital_status_pred |
|-----------|-------------------|
| C3L-03395 | 1 |
| C3N-00709 | 0 |
| C3N-03039 | 0 |
| C3N-03069 | 0 |

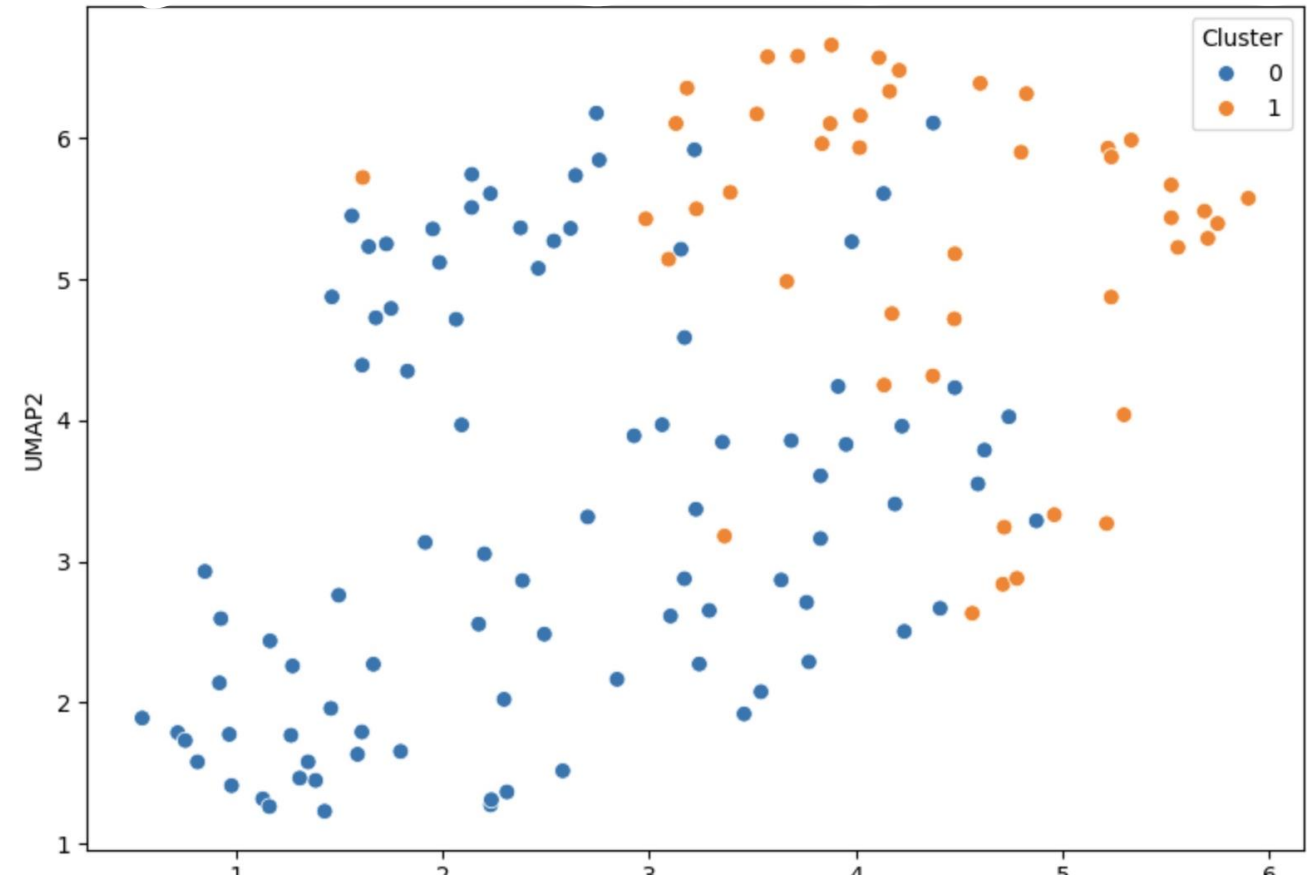
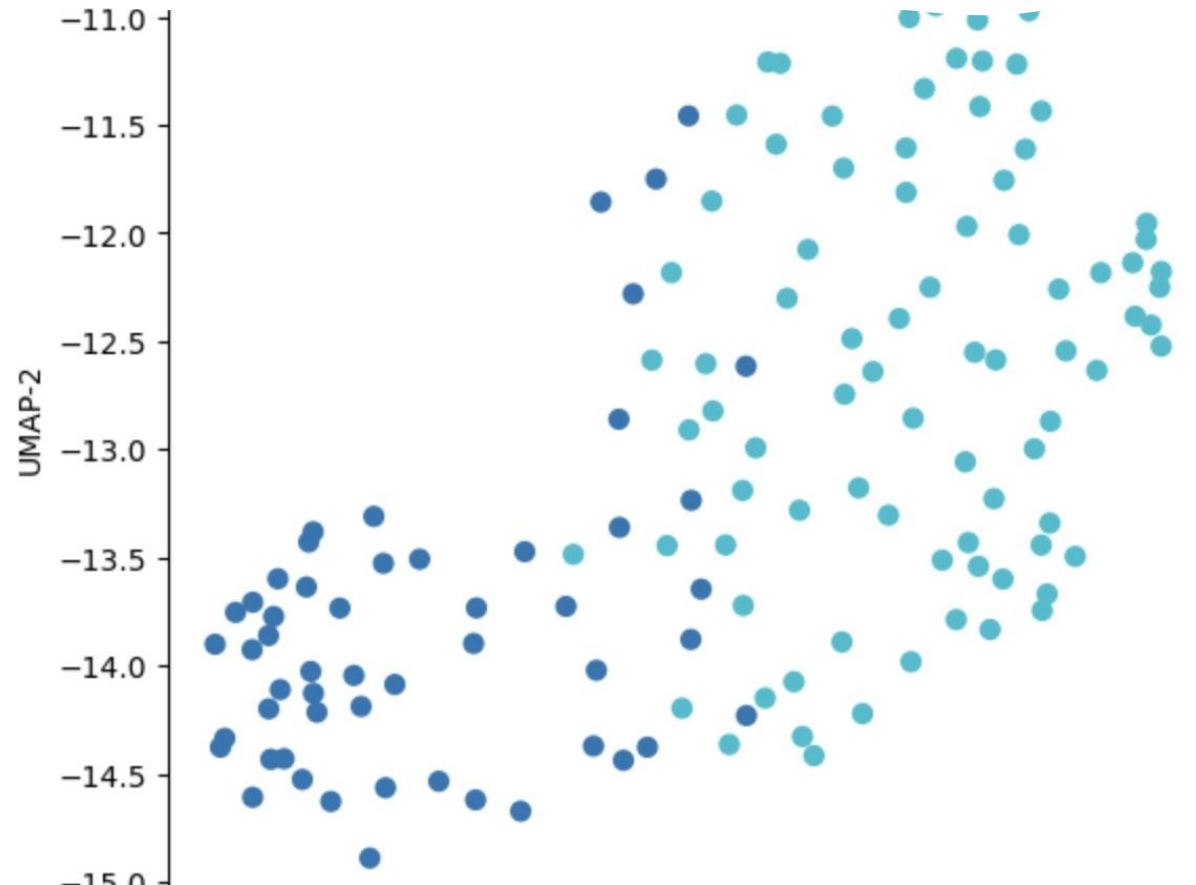
(0 is alive, 1 is deceased)

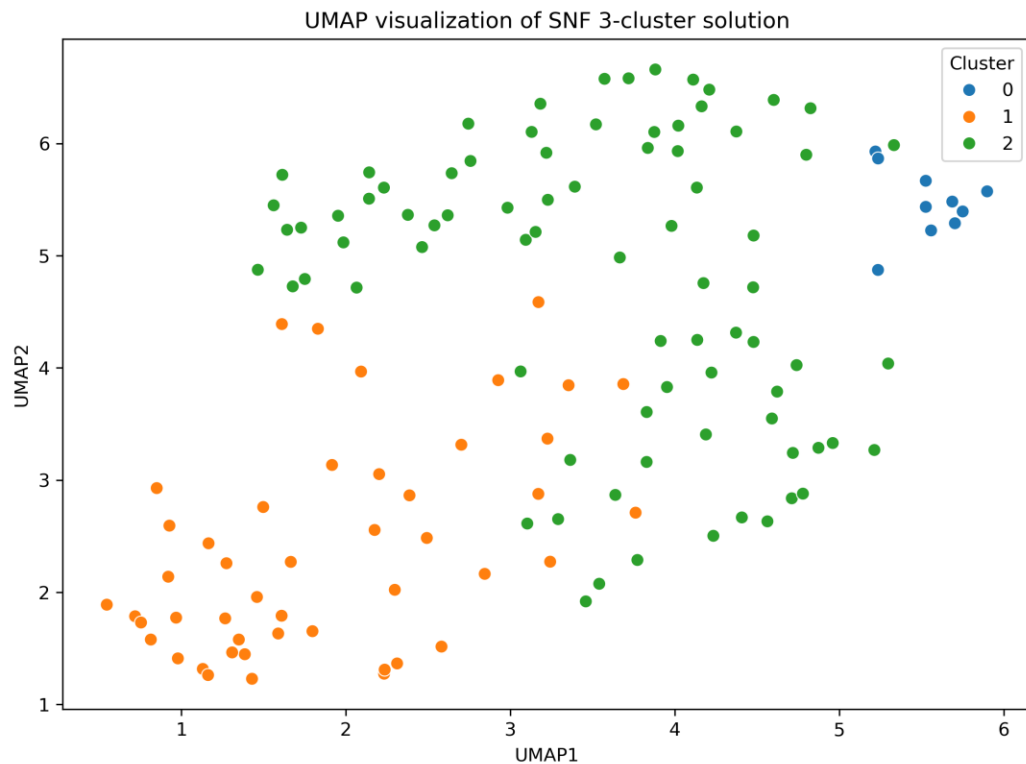


Second step:

The SCNA_gene and mRNA datasets were integrated using SNF in Python and subsequently clustered. The resulting cluster labels were added to the clinical data to perform the same classification, predicting the patients' vital status, and the results were compared with the previous classification outcomes.

- The SCNA_gene dataset and the mRNA dataset were integrated using SNF in Python. After fusion, the data were transformed into a sample-by-sample similarity matrix and subsequently clustered. We initially used the parameters $K = 20$ and $\mu = 0.5$, but this resulted in a low silhouette score. After optimizing the parameters, the clustering performance improved

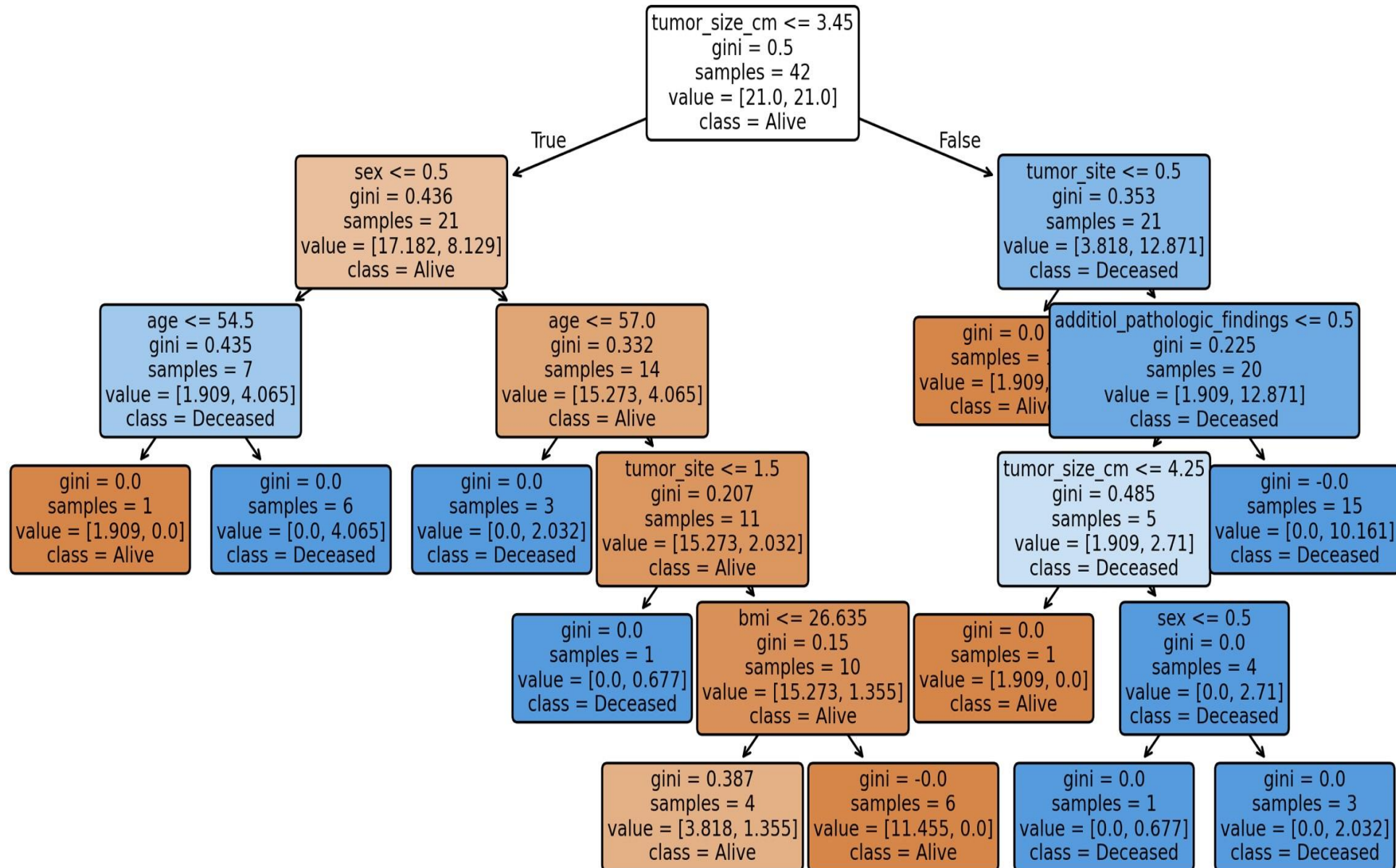




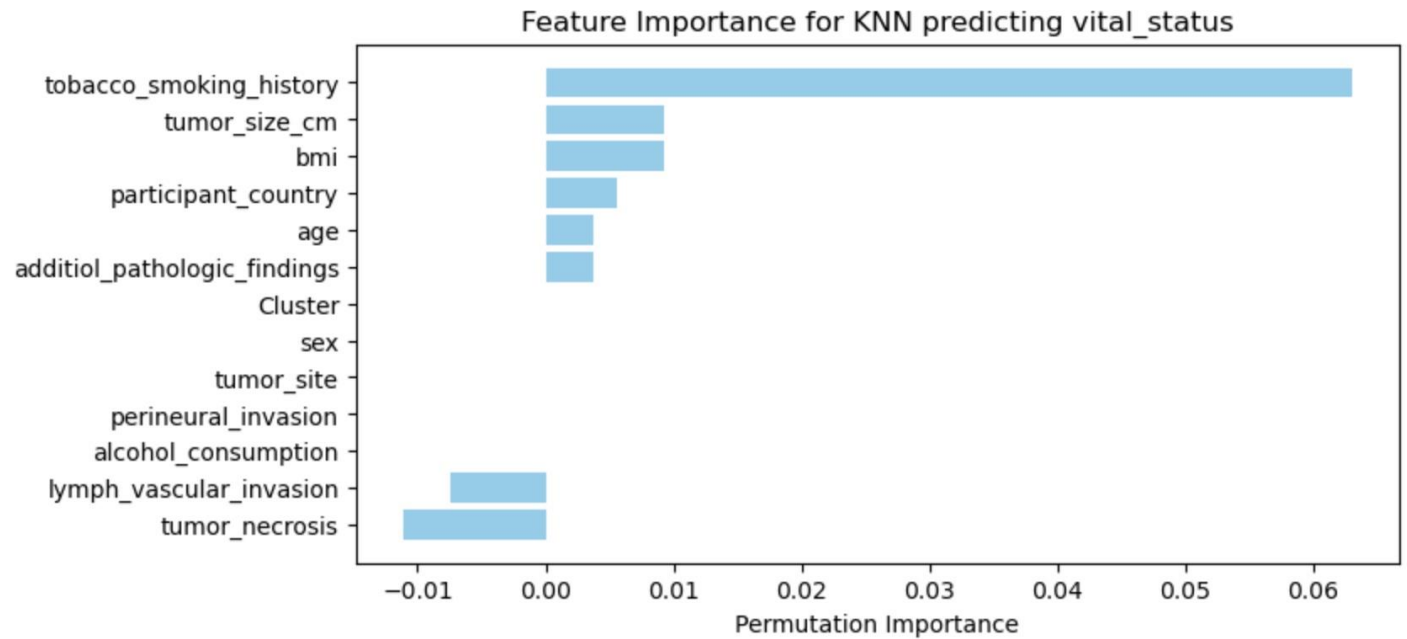
the number of clusters also affected the silhouette score, with the best score achieved when using three clusters.

clusters were still visually distinguishable even when the silhouette score was relatively low.

Decision Tree for Predicting Vital Status



KNN : integrated data



When the clusters were added to the classification, the results changed noticeably. The feature importance in the KNN model shifted substantially, with smoking and BMI remaining the dominant features

Predictet vital status

case_id predicted_vital_status predicted_prob_deceased

C3L-03395 1 0.8

C3N-00709 1 0.6

C3N-03039 1 1.0

C3N-03069 1 1.0

(0 is alive, 1 is deceased)

