

# Prediction of Lung Cancer Stage

ANJA PETKOVIĆ SV22/2020, BRANISLAV STOJKOVIĆ SV64/2020

## 1. MOTIVATION

---

Since we dealt with the same topic in another subject using rule-based systems, we were motivated to explore it further through machine learning. During the development of that project, we consulted with medical experts specializing in lung cancer. From their insights, we realized that not every patient can be diagnosed using a linear set of rules. Each patient requires special consideration and analysis of their findings. Therefore, we found that machine learning could be used to create such a system.

## 2. RESEARCH QUESTIONS

---

While determining the topic and dataset for the project, we encountered various datasets. Some had fewer rows and columns, while others were more extensive. We aimed for a dataset more complex than those used in exercises, with as many rows and columns as possible. After considering these criteria, we found a suitable dataset available at the following link: [Kaggle Lung Cancer Prediction Dataset](#).

The dataset includes the following attributes:

- **Patient\_ID**: Patient identification number.
- **Age**: Age of the patient.
- **Gender**: Gender of the patient.
- **Smoking\_History**: Smoking history of the patient (e.g., Current Smoker, Never Smoked, Former Smoker).
- **Tumor\_Size\_mm**: Tumor size in millimeters.
- **Tumor\_Location**: Tumor location (e.g., Lower Lobe).
- **Stage**: Cancer stage (e.g., Stage I, Stage II, Stage III).
- **Treatment**: Treatment received by the patient (e.g., Surgery, Radiation Therapy, Chemotherapy).
- **Survival\_Months**: Number of survival months after diagnosis.
- **Ethnicity**: Ethnicity of the patient.
- **Insurance\_Type**: Type of insurance (e.g., Medicare, Private, Other).
- **Family\_History**: Family history of cancer (Yes/No).
- **Comorbidity\_Diabetes**: Presence of diabetes (Yes/No).
- **Comorbidity\_Hypertension**: Presence of hypertension (Yes/No).
- **Comorbidity\_Heart\_Disease**: Presence of heart disease (Yes/No).
- **Comorbidity\_Chronic\_Lung\_Disease**: Presence of chronic lung disease (Yes/No).
- **Comorbidity\_Kidney\_Disease**: Presence of kidney disease (Yes/No).
- **Comorbidity\_Autoimmune\_Disease**: Presence of autoimmune diseases (Yes/No).
- **Comorbidity\_Other**: Presence of other comorbidities (Yes/No).
- **Performance\_Status**: Performance status of the patient (scale from 0 to 4).

- **Blood\_Pressure\_Systolic:** Systolic blood pressure.
- **Blood\_Pressure\_Diastolic:** Diastolic blood pressure.
- **Blood\_Pressure\_Pulse:** Pulse.
- **Hemoglobin\_Level:** Hemoglobin level.
- **White\_Blood\_Cell\_Count:** White blood cell count.
- **Platelet\_Count:** Platelet count.
- **Albumin\_Level:** Albumin level.
- **Alkaline\_Phosphatase\_Level:** Alkaline phosphatase level.
- **Alanine\_Aminotransferase\_Level:** Alanine aminotransferase level.
- **Aspartate\_Aminotransferase\_Level:** Aspartate aminotransferase level.
- **Creatinine\_Level:** Creatinine level.
- **LDH\_Level:** Lactate dehydrogenase level.
- **Calcium\_Level:** Calcium level.
- **Phosphorus\_Level:** Phosphorus level.
- **Glucose\_Level:** Glucose level.
- **Potassium\_Level:** Potassium level.
- **Sodium\_Level:** Sodium level.
- **Smoking\_Pack\_Years:** Number of smoking years (expressed in pack-years).

### 3. RELATED WORK

---

Based on this dataset, it is possible to predict several columns, such as the cancer stage, the number of months the patient will survive, and the type of therapy. Most solutions focus on predicting the cancer stage and the number of months the patient will survive. The preprocessing steps for these classifications typically include removing irrelevant features, encoding features, and various types of normalization. After data preparation, models are built. The models used include SVC, RandomForestClassifier, AdaBoostClassifier, DecisionTreeClassifier, KNeighborsClassifier, XGBClassifier, CatBoostClassifier, LGBMClassifier, StackingClassifier, RidgeClassifier, Perceptron, and ExtraTreesClassifier.

### 4. METHODOLOGY

---

To address the problem of predicting the stage of lung cancer, we followed a systematic approach involving data preprocessing, feature engineering, model training, and evaluation. Here is a detailed explanation of the methodology:

#### 4.1 DATA PREPROCESSING

- **Data loading:** The dataset was loaded using pandas.
- **Feature selection:** We dropped irrelevant features that were not contributing to the stage prediction, such as **Survival\_Months**, **Blood\_Pressure\_Systolic**, **Blood\_Pressure\_Diastolic**, **Age**, **Blood\_Pressure\_Pulse**, **Performance\_Status**, **Gender**, **Ethnicity**, **Insurance\_Type** and **Patient\_ID**. Features were dropped based on heatmap and classification report.
- **Splitting data:** The dataset was split into training and testing sets using an 80-20 split.

## 4.2 FEATURE ENGINEERING

- **Categorical features:** Included Smoking\_History, Tumor\_Location, Treatment, Ethnicity, Insurance\_Type, and Family\_History.
- **Numerical features:** Included Tumor\_Size\_mm, Hemoglobin\_Level, White\_Blood\_Cell\_Count, Platelet\_Count, Albumin\_Level, Alkaline\_Phosphatase\_Level, Alanine\_Aminotransferase\_Level, Aspartate\_Aminotransferase\_Level, Creatinine\_Level, LDH\_Level, Calcium\_Level, Phosphorus\_Level, Glucose\_Level, Potassium\_Level, Sodium\_Level, and Smoking\_Pack\_Years.

## 4.3 DATA TRANSFORMATION PIPELINES

- **Numerical transformer:** Applied median imputation followed by standard scaling.
- **Categorical transformer:** Applied most frequent imputation followed by one-hot encoding.

## 4.4 MODEL TRAINING

- **Gradient Boosting Classifier:**
  - **Training:** A Gradient Boosting Classifier was trained using the preprocessed data.
  - **Evaluation:** The model's performance was assessed using macro and micro F1 scores, precision and recall.
- **Stacking Classifier:**
  - **Training:** A stacking classifier combining SVM, KNN, and Gaussian Naive Bayes with Logistic Regression as the final estimator was trained.
  - **Evaluation:** The model's performance was assessed using macro and micro F1 scores, precision and recall.

## 4.5 MODEL EVALUATION

- **Evaluation metrics:** macro and micro F1 score, precision, recall.
- **Hyperparameter optimization:** performed using Random Search algorithm.

# 5. DISCUSSION

---

For testing the results, we used the test dataset, which is 20% of the initial dataset. As performance metrics, we used the Micro and Macro F1 scores, as well as Precision and Recall metrics. Parameter optimization was carried out using the RandomSearch algorithm. Table 1 shows the metric values for both classifiers. Picture 1 represents the classification report for Gradient Boosting classifier. Picture 2 represents the classification report for Stacking classifier.

Table 1 Results

CLASSIFIER	MICRO F1	MACRO F1	PRECISION	RECALL
BOOSTING	0.254	0.250	0.257	0.256
STACKING	0.251	0.207	0.44	0.256

Boosting - Classification Report:				
	precision	recall	f1-score	support
Stage I	0.25	0.23	0.24	1170
Stage II	0.28	0.16	0.21	1233
Stage III	0.27	0.26	0.26	1175
Stage IV	0.24	0.37	0.29	1154
accuracy			0.25	4732
macro avg	0.26	0.26	0.25	4732
weighted avg	0.26	0.25	0.25	4732

Picture 1 Classification report for boosting model

Stacking - Classification Report:				
	precision	recall	f1-score	support
Stage I	0.27	0.20	0.23	1170
Stage II	1.00	0.00	0.00	1233
Stage III	0.25	0.27	0.26	1175
Stage IV	0.24	0.56	0.34	1154
accuracy			0.25	4732
macro avg	0.44	0.26	0.21	4732
weighted avg	0.45	0.25	0.20	4732

Picture 2 Classification report for stacking model

#### KEY OBSERVATIONS FROM THE RESULTS:

- **Gradient Boosting Classifier:** This model demonstrated relatively balanced performance across different stages, with no single stage showing extreme values in terms of precision or recall. This suggests that the model can consistently identify instances of each stage to some extent.
- **Stacking Classifier:** This model showed a high precision but very low recall for Stage II, meaning it predicted Stage II very accurately when it did predict it, but it rarely predicted Stage II. This imbalance indicates that while the model has the potential for accurate predictions, it may require further tuning or balancing to improve recall and overall performance.

## 6. REFERENCES

---

[HTTPS://WWW.KAGGLE.COM/DATASETS/RASHADRMAMMADOV/LUNG-CANCER-PREDICTION](https://www.kaggle.com/datasets/rashadrmammadov/lung-cancer-prediction)