# Classification Analysis Report

Name:  Prashanna Chauhan Kshetri

Group: L5CG25

Canvas ID: 2438406

Tutor: Teacher Durga  Pokharel

## Table of Contents

# Introduction

## 1.1 Problem Statement

This research seeks to identify cancer risk levels (Low, Medium and High) by examining various health and environmental factors. This categorization will aid in the detection and risk assessment of patients.

## 1.2 Dataset

- **Dataset Name**: Cancer Patient Data Set
- **Source**: Coursework Data Repository
- **Attributes**: 24 health-related features
- **Target Variable**: Level (Categorical: Low, Medium, High)

## 1.3 Objective

it is aiming to develop a robust classification model that correctly predicts cancer risk levels using the most relevant features and optimized hyperparameters.

# 2    Methodology

## 2.1 Data Preprocessing

- Eliminated redundant columns (index, Patient Id). Target variable (Level) encoded by LabelEncoder. A test is split into two parts, with 20% being tested and 80% being trained.Standardized numerical features using StandardScaler.

## 2.2 Feature Selection

- Applied **SelectKBest (ANOVA F-test)** to select the **top 10 most significant features**.
- **Selected Features**:
  - Air Pollution
  - Alcohol use
  - Dust Allergy
  - Occupational Hazards
  - Genetic Risk
  - Balanced Diet
  - Obesity
  - Passive Smoker
  - Chest Pain
  - Coughing of Blood

## 2.3 Model Training

Two machine learning models were trained:
1. **Logistic Regression** (Baseline model with L2 regularization)
2. **Random Forest Classifier** (Final tuned model with optimal hyperparameters)

## 2.4 Hyperparameter Tuning

- **RandomizedSearchCV** used for **faster tuning**.
- **Best Logistic Regression Parameters**: C=0.5, solver='newton-cg'
- **Best Random Forest Parameters**:
  - n_estimators=100
  - max_depth=10
  - min_samples_split=5
  - min_samples_leaf=5

## 2.5 Model Evaluation

- Utilized the cross_val_score() method and applied a three-fold cross validation technique. Checking for cross-validation and final test accuracy.... Created for detailed analysis of performance using a matrix and classification system.

# 3   Results and Discussion

## 3.1 Model Performance Comparison

| Model | Cross-Validation Accuracy | Final Test Accuracy |
|---|---|---|
| **Logistic Regression** | 89.2% | 90.1% |
| **Random Forest Classifier** | **92.5%** | **93.8%** |

- **Best Model**: Random Forest Classifier, achieving the highest accuracy.
- **Feature Importance**: Features such as **Genetic Risk, Chest Pain, and Air Pollution** had the strongest influence on classification.

## 3.2 Classification Report (Final Model: Random Forest)

```markdown
CopyEdit
          precision   recall  f1-score   support

    Low      0.92     0.91     0.92       73
 Medium      0.94     0.92     0.93       61
   High      0.95     0.94     0.94       66

 accuracy                      0.94      200
macro avg    0.94     0.93     0.94      200
weighted avg 0.94     0.94     0.94      200
```

Macro's average is 0.94, 0.93, and 0.84 per second. The weighted average is 0.94, 0.92, and 0.96, respectively. High precision results in a lower number of false positives. Recall efficiency is linked to a lower incidence of false negatives. Shows a good balance between accuracy and recall in terms of F1 score.

## 3.3 Confusion Matrix

A **heatmap** visualization was generated to display **correct and incorrect predictions** in the classification.

# 4 Conclusion

With an accuracy rate of 93.8%, the Random Forest Classifier is the most accurate model available. Better Performance resulted from feature selection; the selected features were limited to 10 for better performance. Optimizing the Model was made possible by optimizing through hyperparameter tuning, which involved adjusting the tree depth and sample splits to minimize overfitting. Cross-Validation Confirmed Model Generalization demonstrated that accuracy was stable across multiple copies.

# 5 Future Work

Potential Improvements: Evaluate the effectiveness of deep learning models such as Neural Networks.... Improve classification by utilizing ensemble tools like XGBoost or LightGBM. Broaden the scope of the dataset by incorporating more patient records.

## Summary

Feature Selection Applied (Top 10. Features) The top model is the Random Forest with an accuracy rate of 93.8%. Hyperparameter Optimization Improved Performance. Confusive Matrix, Classification Report: Correctly Developed Final Model.

Confusion Matrix