# Milestone 1 Report

Data Collection, Preprocessing, and Exploratory Data Analysis (EDA)

**Project Objective**

The goal is to develop a heart disease risk prediction model and interactive dashboard based on stress levels, sleep duration, and gym exercise patterns. By integrating multiple datasets, this model aims to identify key factors contributing to heart disease risk and provide actionable insights.

**Datasets Used**

1. Sleep Health and Lifestyle Dataset:

   - Source: Kaggle (https://www.kaggle.com/datasets/uom190346a/sleep-health-and-lifestyle-dataset)
   - Attributes: Person ID, Gender, Age, Occupation, Sleep Duration, Quality of Sleep, Physical Activity Level, Stress Level, BMI Category, Blood Pressure, Heart Rate, Daily Steps, Sleep Disorder

2. Heart Disease Dataset

   - Source: Kaggle (https://www.kaggle.com/datasets/johnsmith88/heart-disease-dataset)
   - Attributes: Age, gender, chest pain type, blood pressure, cholesterol, fasting blood sugar, ECG results, maximum heart rate, heart disease diagnosis.


3. Gym Members Exercise Tracking Dataset
   - Source: Kaggle (https://www.kaggle.com/datasets/valakhorasani/gym-members-exercise-dataset)
   - Attributes: Age, Gender, Workout frequency, calories burned, BMI, hydration levels, heart rate, daily steps.

Common Attributes are Age and Gender and Heart Rate

**Tech Stack**

Programming Language

- Python

Libraries & Frameworks

- **Data Manipulation**: Pandas, NumPy

- **Visualization**: Matplotlib, Seaborn

**Data Preprocessing**

Steps taken for data preprocessing include:

- Handling missing values through median imputation.

- Splitting 'Blood Pressure' into 'Systolic_BP' and 'Diastolic_BP'.

- Standardizing column names for consistency.

- Encoding categorical variables such as gender and BMI categories.

- Detecting and handling outliers using the IQR method.

- Merging datasets on age and gender to form a unified dataset.

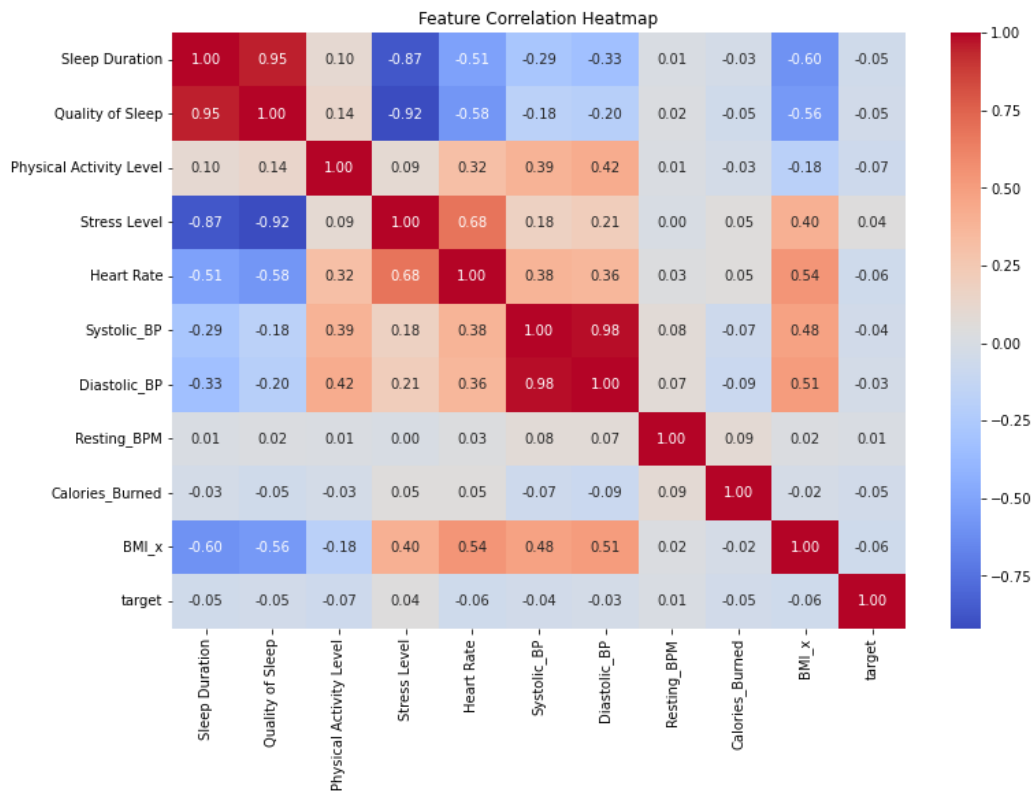- Removing duplicate rows to ensure data integrity.

**Exploratory Data Analysis (EDA)**

The following analyses were performed:

- Descriptive statistics: Summary statistics computed for all numerical variables.

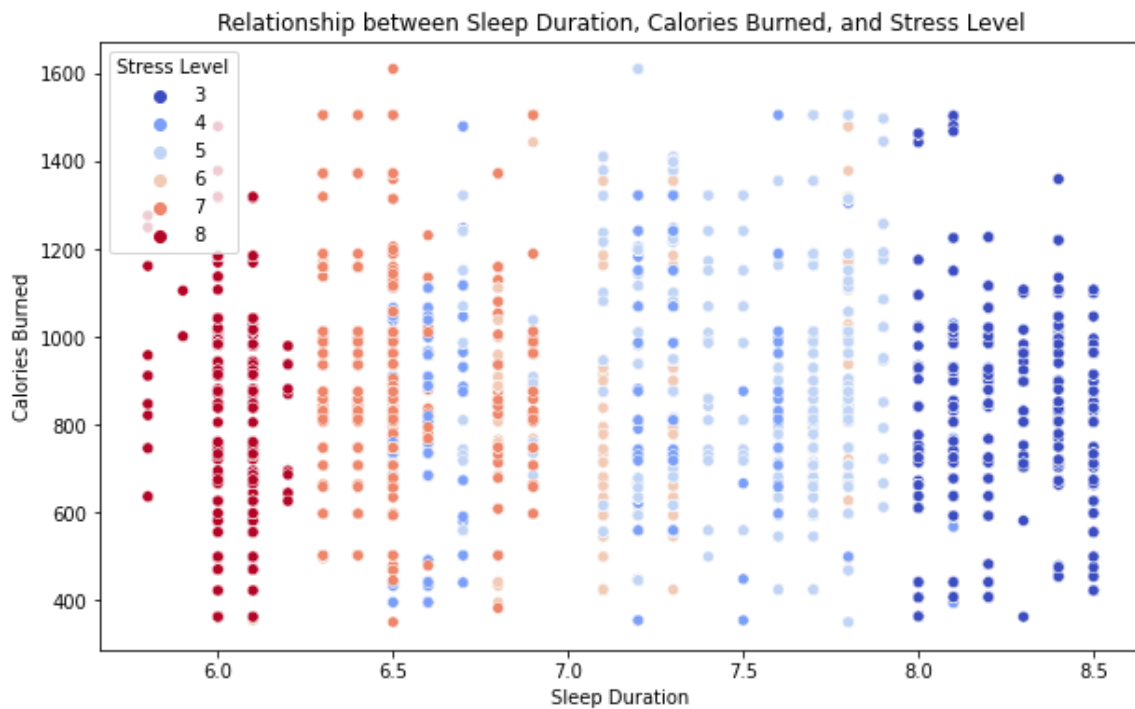- shape of balanced dataset is 6386 rows and 39 columns

**Important Visualizations:**

 - Correlation Heatmap of merged dataset to identify feature relationships. This heatmap tells that **stress levels, sleep duration, and physical activity are key factors** in determining heart health.
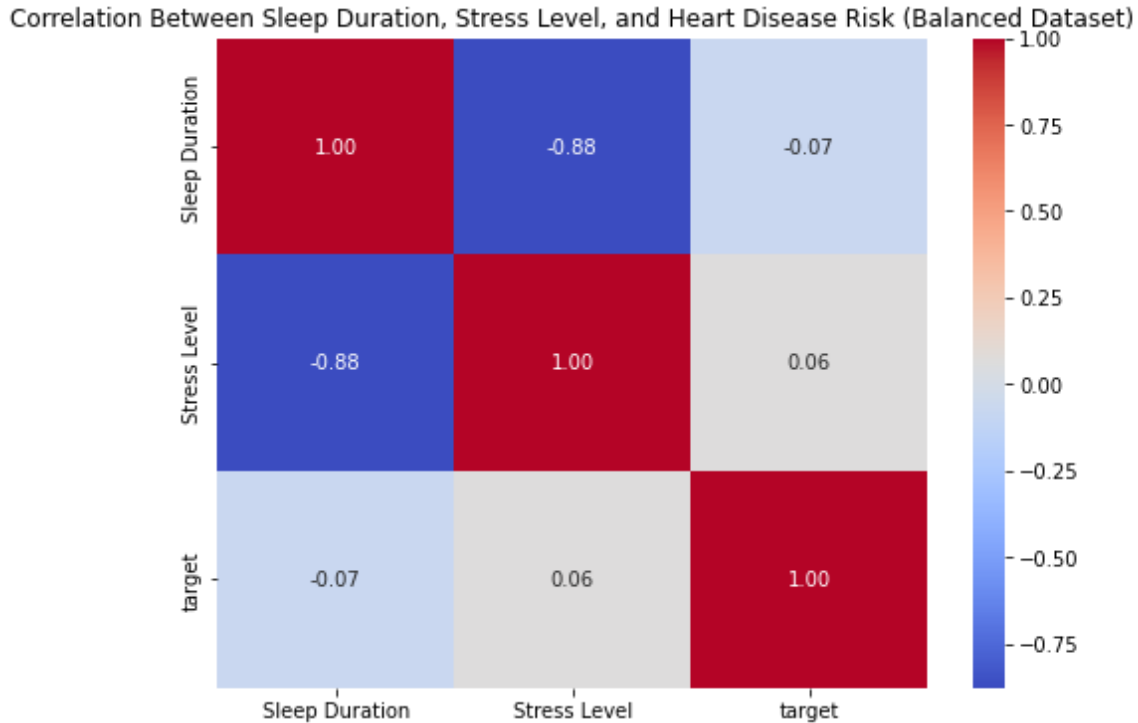


Feature Correlation Heatmap

- Scatter plot to analyze relationships between sleep duration, calories burned, and stress level.

- This scatter plot visualizes stress levels across different individuals. Higher stress levels (darker red) tend to cluster on the left side, while lower stress levels (blue) are more evenly spread across the distribution.

- The separation between high-stress and low-stress groups suggests a strong impact of lifestyle habits on stress regulation. High-stress individuals may require more intervention strategies to improve heart health.



Relationship between Sleep Duration, Calories Burned, and Stress Level

- Correlation Heatmap between "Sleep Duration, Stress Level, and Heart Disease Risk (Balanced Dataset) interprets strong relationships between stress levels, sleep duration, and heart rate. **High stress correlates with increased heart rate and reduced sleep duration.**



Correlation Between Sleep Duration, Stress Level, and Heart Disease Risk (Balanced Dataset)

- Class imbalance handling: The dataset was balanced using undersampling to improve model performance.

Balanced Dataset details:
- No missing values after preprocessing.
- Outliers handled appropriately.
- Balanced dataset for improved model performance in the next stages.