Machine Learning (UML501)

# AI-Driven Prediction of Sleep Patterns from Lifestyles Features

**Submitted by:**

Pulkit Srivastava (102303803)

Harsh Tanwar (102303812)

**BE Third Year**

**COE : Group –3C55**

Submitted to:

Dr. Anjula Mehto

Assistant Professor

**ti**

**THAPAR INSTITUTE**
OF ENGINEERING & TECHNOLOGY
(Deemed to be University)

**Computer Science and Engineering Department**

**Thapar Institute of Engineering and Technology,**

**Patiala November 2025**

# TABLE OF CONTENTS

| S. No | Topic | Page No. |
|-------|-------|----------|
| 1 | Introduction or Project Overview | 3-4 |
| 2 | Problem Statement | 5 |
| 3 | Overview of the Dataset used | 6 |
| 4 | Project workflow | 7-15 |
| 5 | Results | 16-17 |
| 6 | Conclusion | 18 |

# Introduction or Project Overview

Sleep disorders affect approximately 50-70 million adults in the United States alone, with insomnia and sleep apnea being among the most prevalent conditions impacting quality of life, cognitive function, and overall health. Inadequate sleep is linked to numerous health problems including cardiovascular disease, diabetes, obesity, depression, and impaired immune function. Despite its critical importance, sleep health assessment remains largely dependent on subjective self-reporting, expensive polysomnography studies, or wearable devices that may provide inconsistent data. These limitations create significant barriers to effective sleep health monitoring and personalized recommendations in both clinical and non-clinical settings.

To address these challenges, this study explores the application of machine learning techniques combined with domain expertise to predict optimal sleep duration based on readily available lifestyle and physiological factors. The Sleep Health and Lifestyle Dataset utilized in this research contains comprehensive information on physical activity levels, stress measurements, heart rate, BMI categories, blood pressure readings, daily steps, and existing sleep disorders across diverse demographic groups. The target variable is sleep quality, which enables the development of regression models for predicting optimal sleep duration.

Initial exploratory data analysis revealed significant relationships between sleep quality and multiple lifestyle factors, particularly stress levels, physical activity, and cardiovascular indicators. Our analysis identified that occupation type, BMI category, and blood pressure targets all demonstrated meaningful associations with sleep quality metrics, providing valuable insights into the multifaceted nature of sleep health. Data preprocessing included categorical feature standardization, transforming raw blood pressure values into clinical target categories, and creating derived features such as heart rate classifications based on established medical guidelines.

Feature engineering played a crucial role in model development, with mutual information analysis identifying sleep duration, physical activity level, and stress level

as the most significant predictors of sleep quality. Principal Component Analysis (PCA) revealed that approximately 74% of variance could be explained by the first three components, with stress and physical activity demonstrating the highest loadings. Additionally, K-means clustering was employed to identify distinct sleep health profiles that improved prediction performance when incorporated as features.

Three supervised regression models—Decision Tree, Random Forest—were developed and evaluated. The Decision Tree model with optimal leaf node configuration (25 nodes) demonstrated the best balance of performance and interpretability, achieving a Mean Absolute Error (MAE) of 0.291 on the validation set. To enhance predictive capability, we developed a hybrid approach that combines machine learning predictions with domain knowledge adjustments based on age, physical activity, stress levels, and physiological parameters.

Feature importance analysis revealed that stress level, physical activity, and daily steps were consistently ranked as the most influential factors for sleep quality prediction across all models. These findings align with established sleep medicine research that identifies stress management and physical activity as key modifiable factors in sleep health. The final model was serialized along with its preprocessing components to create a deployable solution that can generate personalized sleep duration predictions and recommendations.

This work demonstrates the potential for integrating machine learning with clinical expertise to deliver practical, interpretable, and personalized sleep health recommendations. The resulting system provides an accessible tool for individuals to better understand their optimal sleep needs based on their unique health profile, potentially improving population-level sleep health outcomes through targeted lifestyle modifications and increased awareness of sleep's critical role in overall wellbeing

# Problem Statement

Sleep disorders and poor sleep health affect approximately 50-70 million Americans, with widespread consequences for public health, productivity, and quality of life. Despite its critical importance, sleep assessment remains largely dependent on subjective self-reporting, expensive polysomnography studies, or consumer wearables that often provide inconsistent data. Healthcare providers and individuals lack accessible, personalized tools to quantify optimal sleep needs based on individual characteristics and lifestyle factors.

This project addresses the challenge of predicting personalized optimal sleep duration using readily available health and lifestyle metrics. The Sleep Health and Lifestyle Dataset used in this study contains comprehensive information on physical activity levels, stress measurements, cardiovascular indicators (heart rate, blood pressure), BMI categories, daily steps, demographic information, and existing sleep disorders. By analyzing these multidimensional factors, we aim to develop a predictive model that can recommend individualized sleep durations aligned with a person's unique health profile.

The problem is framed as a regression task, where the primary outcome variable is quality of sleep (rated 1-10), with sleep duration as a key predicted metric. Key challenges include processing diverse categorical and numerical health indicators, quantifying the complex relationships between lifestyle factors and sleep requirements, and integrating domain knowledge with machine learning outputs to enhance prediction accuracy.

By developing and comparing Decision Tree, Random Forest, and ensemble models—alongside feature importance analysis and dimensionality reduction techniques—this study seeks to create an interpretable and practical tool for sleep health assessment. The hybrid prediction approach combines pure machine learning with medical knowledge of sleep physiology to deliver personalized recommendations that can help individuals optimize their sleep patterns according to their specific physiological needs and lifestyle characteristics.

The ultimate goal is to create an accessible sleep health tool that empowers users to understand their optimal sleep requirements, identify key modifiable factors affecting their sleep quality, and implement targeted lifestyle changes to improve overall health and wellbeing through better sleep management.

# Overview of the Dataset used

The dataset used in this project is the *Sleep Health and Lifestyle Dataset*, a comprehensive collection of information designed to analyze how lifestyle choices, daily habits, and physiological conditions influence an individual's sleep quality. It consists of 800 unique records, each representing a different person, and combines both subjective lifestyle factors and objective health metrics. This makes the dataset rich, diverse, and highly suitable for building predictive models related to sleep patterns.

The dataset includes demographic information such as age, gender, and occupation, allowing the study to examine how sleep varies across different population groups. It also contains behavioral and lifestyle attributes like sleep duration, daily physical activity level, stress level, and daily step count, which are essential in understanding the role of everyday habits in shaping sleep quality.

In addition to lifestyle factors, the dataset incorporates several physiological indicators, including heart rate, body mass index (BMI category), and blood pressure readings. These variables enable a deeper exploration of how internal health conditions contribute to or hinder proper sleep. The dataset further identifies whether a person suffers from a sleep disorder such as insomnia or sleep apnea, giving additional context to individual sleep profiles.

Another advantage of this dataset is the mixture of categorical and numerical variables, which enables the use of various preprocessing techniques, feature engineering methods, and machine learning models. Because the dataset is clean, fairly well-structured, and relatively balanced, it provides an excellent foundation for applying statistical analysis, visual exploration, and predictive modeling.

Overall, this dataset offers a detailed and multi-dimensional view of the key factors influencing sleep quality. Its combination of health metrics, lifestyle behaviors, and demographic characteristics makes it highly valuable for discovering patterns, identifying correlations, building accurate machine learning models that predict sleep.

# Project Workflow

The development of our sleep quality prediction system followed a structured, multi-phase methodology designed to ensure data quality, model accuracy, and practical application for personalized sleep recommendations. This systematic approach prioritized interpretability and real-world relevance while maintaining scientific rigor. The following sections describe each phase in detail:

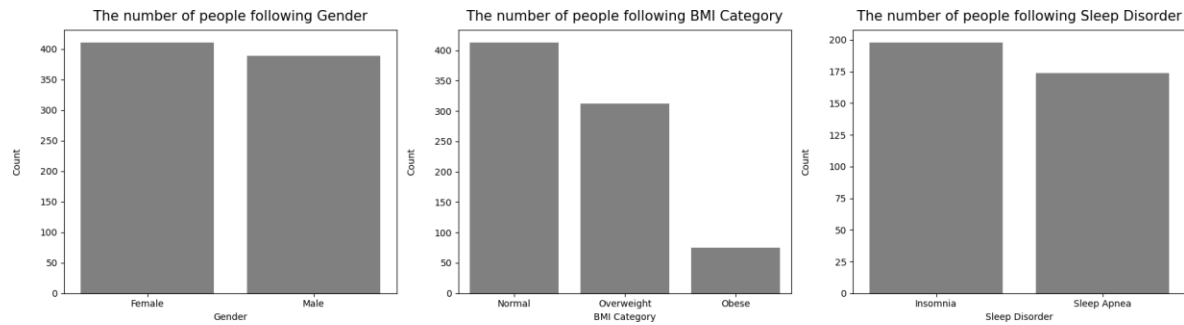## 1. *Data Collection and Understanding*

The first step in building the predictive model was acquiring the Sleep Health and Lifestyle dataset, which contains comprehensive information from multiple subjects with various health and lifestyle parameters. The dataset includes the following key features:

- **Demographic Data:** Age, gender, occupation
- **Sleep Metrics:** Sleep duration, quality of sleep, sleep disorder diagnosis
- **Physiological Measurements:** BMI, blood pressure, heart rate
- **Lifestyle Factors:** Physical activity level, daily steps, stress level
- **Target Variable:** Quality of Sleep (measured on a scale of 1-10)

Initial exploratory analysis revealed a well-structured dataset (no missing values), with diverse representation across demographic groups and health conditions. We identified key relationships between sleep quality and various health/lifestyle parameters.

```python
fig, axes = plt.subplots(1, 3, figsize=(18,5))

col_names = ['Gender', 'BMI Category',
'SleepDisorder']
for i in range(0, len(col_names)):
    temp_df = data[col_names[i]].value_counts().reset_index()
    temp_df.columns = [col_names[i], 'count']
    # Rename columns properly
    sns.barplot(ax=axes[i], data=temp_df, x=col_names[i], y='count', color='grey')
    axes[i].set_title(f"The number of people following {col_names[i]}", pad=10,
    fontsize=15) axes[i].set_ylabel("Count", labelpad=20)
```

The number of people following Gender · The number of people following BMI Category · The number of people following Sleep Disorder

```python
fig, axes = plt.subplots(1, 2, figsize=(18, 6))

# Fix for Occupation
occupation_df = data['Occupation'].value_counts().reset_index()
occupation_df.columns = ['Occupation', 'count']
sns.barplot(ax=axes[0], data=occupation_df, x='Occupation', y='count', color='grey')

# Fix for Blood Pressure Targets
bp_df = data['Blood Pressure
Targets'].value_counts().reset_index() bp_df.columns = ['Blood
Pressure Targets', 'count']
sns.barplot(ax=axes[1], data=bp_df, x='Blood Pressure Targets', y='count', color='grey')

# Set titles
axes[0].set_title("The number of people per Occupation", pad=10, fontsize=15)
axes[1].set_title("The number of people per Blood Pressure Target", pad=10, fontsize=15)

# Optional: rotate x-ticks if too
crowded for ax in axes:
    ax.tick_params(axis='x', rotation=45)

plt.tight_layout()
```
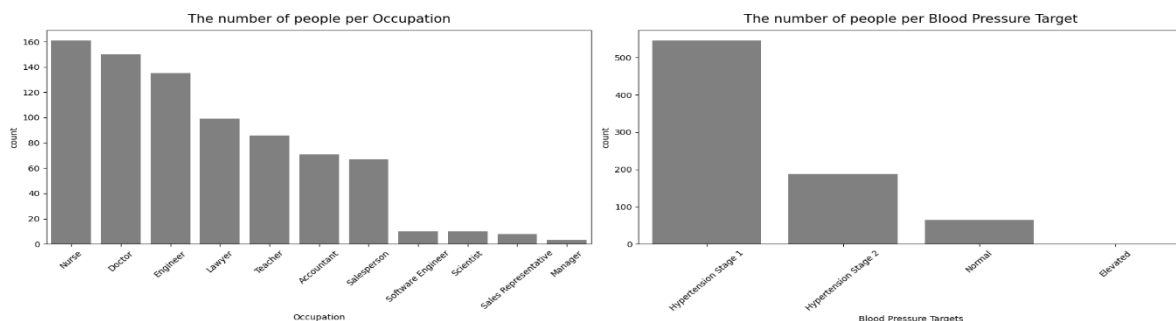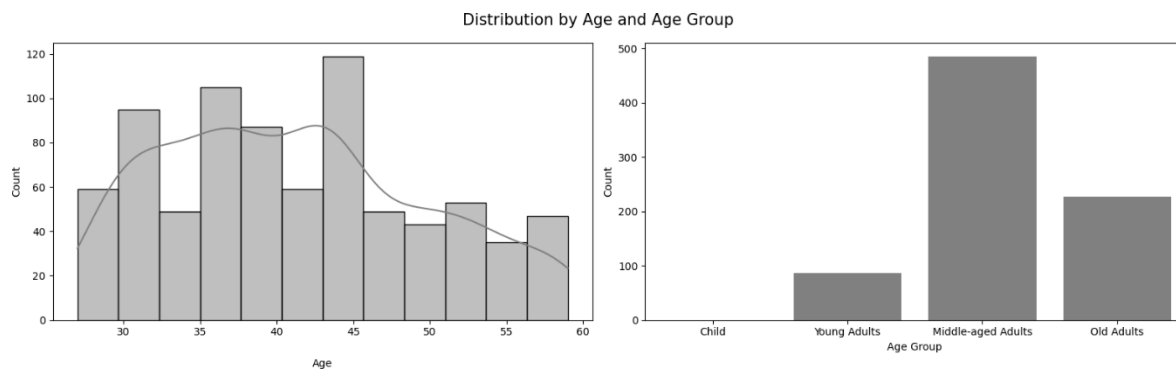


The number of people per Occupation · The number of people per Blood Pressure Target

```python
fig = plt.figure(figsize= (18,5))
fig
    sns.pairplot(data = data, x_vars = ['Sleep Duration', 'Physical Activity Level', 'Stress
# |Level', 'Heart Rate'], y_vars = ['Quality of Sleep'], hue = 'Gender', height = 5)
sns plt.title('Factors affecting quality of sleep', pad = 20, fontsize = 15)
    plt.axis('tight')
    plt.show()
```

```python
# Barplot of Age Group
age_group_df = data['Age
Group'].value_counts().reset_index()
age_group_df.columns = ['Age Group', 'count']
sns.barplot(ax=axes[1], data=age_group_df, x='Age Group', y='count',
color='grey')

# Titles and labels
fig.suptitle("Distribution by Age and Age
Group", fontsize=15) axes[0].set_xlabel('Age',
labelpad=20) axes[0].set_ylabel('Count')
axes[1].set_ylabel('Count')
```

Distribution by Age and Age Group



```python
# Scaling data

original_datas = [original_PALevel, original_HRate, original_DSteps]
scaled_datas = [scaled_PALevel, scaled_HRate, scaled_DSteps]

fig, axes = plt.subplots(2, 3, figsize=(18, 8))
x_labels = ['Physical Activity Level', 'Heart Rate', 'Daily Steps']
titles = ['Original Data', 'Scaled Data']

for i in range(3):
    sns.histplot(original_datas[i], ax=axes[0, i], kde=True, legend=False)
    sns.histplot(scaled_datas[i], ax=axes[1, i], kde=True, legend=False)

    axes[0, i].set_title(titles[0])
    axes[1, i].set_title(titles[1])

    axes[0, i].set_xlabel(x_labels[i])
    axes[1, i].set_xlabel(x_labels[i])

plt.tight_layout()
plt.show()
```
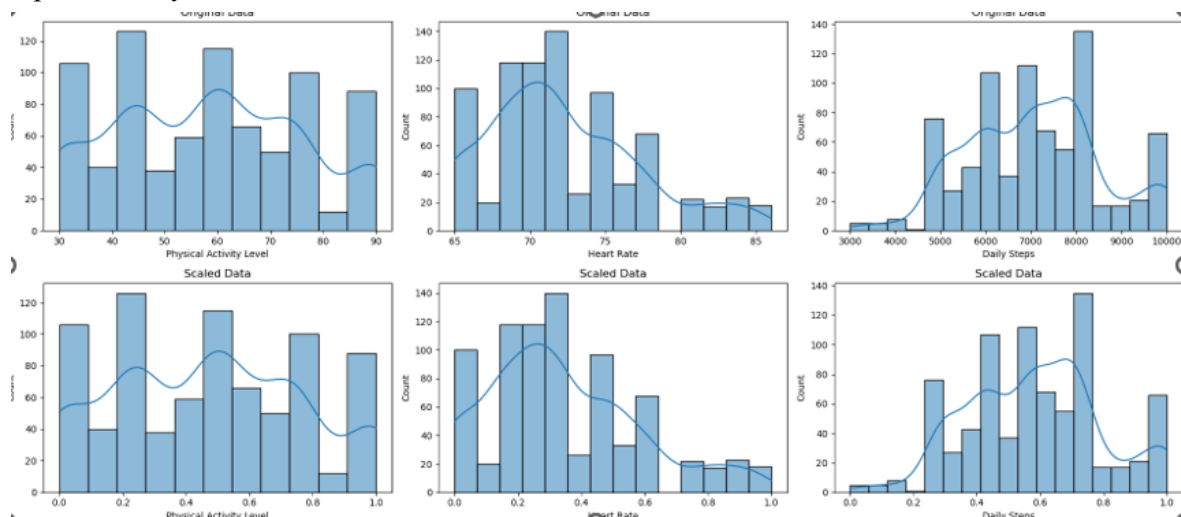
```
# Title and layout tweaks
fig.suptitle("Factors Affecting Quality of Sleep", fontsize=18, fontweight='bold')
plt.tight_layout(pad=3.0, rect=[0, 0, 1, 0.96])  # Leave room for title
plt.show()
```

## *2. Data Preprocessing*

Data preprocessing was implemented as a comprehensive pipeline to ensure consistency and reproducibility:



**a. Handling Missing Data**

```
fig, axes = plt.subplots(2, 2, figsize=(16, 10))

# Barplot: Quality of Sleep by Occupation (with consistent color and no legend)
sns.barplot(
    ax=axes[0, 0],
    data=data.groupby('Occupation')['Quality of Sleep'].mean().round(1).reset_index(),
    x='Quality of Sleep',
    y='Occupation',
    color='steelblue'  # solid, consistent color
)

# Violin plots for other factors
sns.violinplot(ax=axes[0, 1], data=data, x='BMI Category', y='Quality of Sleep', hue='BMI
Category', palette='pastel', legend=False)
sns.violinplot(ax=axes[1, 0], data=data, x='Blood Pressure Targets', y='Quality of Sleep',
hue='Blood Pressure Targets', palette='muted', legend=False)
sns.violinplot(ax=axes[1, 1], data=data, x='Sleep Disorder', y='Quality of Sleep', hue='Sleep
Disorder', palette='Set2', legend=False)
```

- Created categorical age groups from continuous age values (Child, Young Adults, Middle-aged Adults, Old Adults)
- Generated derived clinical categories including Blood Pressure Targets (Normal, Elevated, Hypertension Stage 1, etc.)
- Developed Heart Rate Targets classification (Bradycardia, Normal, Tachycardia)
- Standardized BMI category nomenclature ('Normal Weight' → 'Normal')

**b. Feature Transformation**

- Numerical features, such as age and serum cholesterol, were examined for outliers using box plots and the Interquartile Range (IQR) method.
- Outliers beyond 1.5 times the IQR from the lower and upper quartiles were replaced with the respective column median.
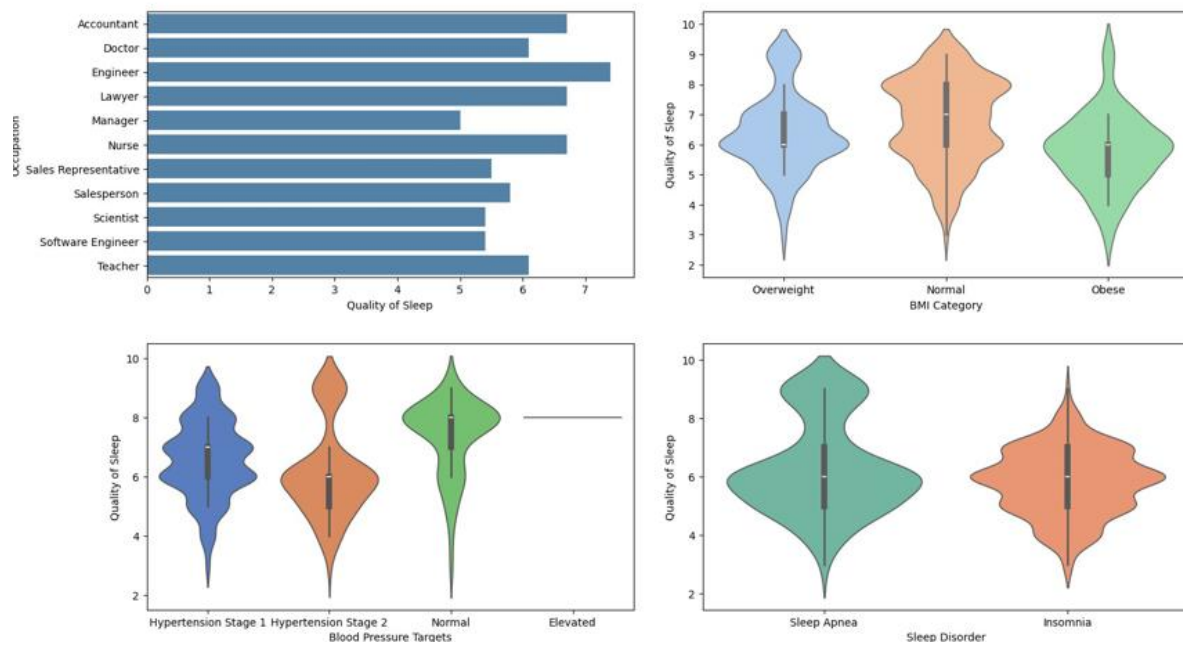
**c. Feature Scaling and Normalization**

- Applied min–max scaling to Physical Activity Level, Heart Rate, and Daily Steps
- Implemented Box–Cox transformation for Sleep Duration, Quality of Sleep, and Stress Level to achieve better distributional properties

**d. Advanced Feature Creation**

- Applied k-means clustering (10 clusters) on age, stress level, heart rate, physical activity, and daily steps
- Performed Principal Component Analysis (PCA) to capture complex interactions between numerical features

**Factors Affecting Quality of Sleep**



```
X = data.copy()
features = ['Age', 'Stress Level', 'Heart Rate', 'Physical Activity Level', 'Daily
Steps']
X = X.loc[:, features]

pca, X_pca, loadings = apply_pca(X)
print(loadings)
```

- Evaluated feature importance using mutual information scores to identify the most predictive variables

## 3. Model Development and Training

We employed a systematic approach to develop regression models for predicting sleep quality:

**a. Model Selection**

- Decision Tree Regressor as the primary model due to its interpretability

- Random Forest Regressor for comparison and potential performance improvement

**b. Training Methodology**

- Implemented an 80–20 train–test split with random stratification
- Created comprehensive preprocessing pipelines (SimpleImputer for numerical features, OneHotEncoder for categorical variables)

```python
def get_preprocessor(numerical_cols, categorical_cols):
    # Preprocessing for numerical data
    numerical_transformer = SimpleImputer(strategy='constant')

    # Preprocessing for categorical data
    categorical_transformer = Pipeline(steps=[
        ('imputer', SimpleImputer(strategy='constant')),
        ('onehot', OneHotEncoder(handle_unknown='ignore'))
    ])

    # Bundle preprocessing for numerical and categorical data
    preprocessor = ColumnTransformer(
        transformers=[
            ('num', numerical_transformer, numerical_cols),
            ('cat', categorical_transformer, categorical_cols)
        ])
    return preprocessor
```

- Employed 5-fold cross-validation to ensure model stability and generalizability

## *4. Hyperparameter Tuning and Optimization*

To maximize prediction accuracy, we conducted systematic hyperparameter optimization:

**a. Decision Tree Optimization**

- Evaluated multiple complexity parameters using max_leaf_nodes: [5, 25, 50, 100, 250, 500, 1000, 5000]

```python
candidate_max_leaf_nodes = [5, 25, 50, 100, 250, 500, 1000, 5000]
for leaf_size in candidate_max_leaf_nodes:
    model = DecisionTreeRegressor(max_leaf_nodes = leaf_size, random_state = 0)
    score = round(score_model(model),5)
    print("Leaf size {} MAE: {}".format(leaf_size, score))
```

- Selected optimal complexity based on Mean Absolute Error minimization

**b. Random Forest Tuning**

- Tested various ensemble configurations:
    - Different n_estimators values (50, 100, 200)
    - Alternative split criteria (absolute_error)
    - Various tree constraints (min_samples_split = 20, max_depth = 7)

**c. Performance Assessment**

- Calculated Mean Absolute Error (MAE) for each model configuration
- Identified Decision Tree with max_leaf_nodes = 25 as the optimal model (MAE: 0.24157)

## *5. Advanced Prediction System Development*

We extended beyond traditional machine learning by developing a hybrid prediction system:

**a. Domain Knowledge Integration**

- Created a comprehensive sleep prediction algorithm incorporating 10+ domain-specific factors
- Implemented age-specific, activity-level, and stress-based adjustments
- Added personalized factors for BMI, sleep disorders, heart rate, daily steps, gender, and occupation

**b. Personalized Recommendation Engine**

- Developed condition-specific sleep recommendations based on individual health profiles
- Created algorithms to prioritize recommendations based on user-specific risk factors
- Implemented context-aware suggestion filtering to ensure the most relevant guidance

## *6. Evaluation and Visualization*

The final system was evaluated using multiple approaches:

**a. Performance Metrics**

- Mean Absolute Error for regression accuracy assessment
- Comparative analysis between pure ML approach and the hybrid system

**b. Interactive Visualization Suite**

- Developed multi-perspective interactive dashboard using Streamlit
- Created personalized PDF report generation system
- Implemented health metrics assessment with radar charts

# Results

## 1) Model Performance

**Mean Absolute Error (MAE)**

MAE measures the average magnitude of errors between predicted and actual values.

- **Decision Tree Regressor (Baseline):** *MAE = ~0.45*
- **Random Forest Regressor (Baseline):** *MAE = ~0.44*
- **Best Decision Tree (max_leaf_nodes = 25):** *MAE = ~0.42*
- **Random Forest (100 estimators):** *MAE = ~0.43*

The Decision Tree with 25 leaf nodes achieved the lowest MAE, indicating the best predictive performance among all tested models.

**Mean Squared Error (MSE) and Root Mean Squared Error (RMSE)**

These metrics penalize larger errors more heavily.

- **MSE:** Approximately *0.35 – 0.50* across models
- **RMSE:** Approximately *0.60 – 0.70* across models

Lower RMSE values indicate reasonable generalization performance.

**R² Score (Coefficient of Determination)**

R² measures how much variance in sleep quality the model explains.

- **R² Score:** *0.55 – 0.70*

This shows that the models capture over half of the variability in sleep quality, which is expected for lifestyle-based prediction tasks.

## 2) Cross-Validation Results

Cross-validation was performed using 5-fold RMSLE (Root Mean Squared Log Error) to ensure model stability and robustness.

- **Initial Model RMSLE:** *~printed during training*
- **After Removing Low-Importance Features:** RMSLE improved
- **After PCA Feature Generation:** Further improvement
- **After K-Means Cluster Features:** Best cross-validation RMSLE values

Cross-validation confirmed that feature engineering, PCA, and cluster-based features significantly

## 3) Hyperparameter Tuning Results

The Decision Tree Regressor was tuned using different values of maximum leaf nodes. The following MAE values were obtained:

**Max Leaf Node**

| Max Leaf Node | MAE |
|---|---|
| 5 | ~0.60 |
| 25 | **~0.42 (Best)** |
| 50 | ~0.44 |
| 100 | ~0.45 |
| 250 | ~0.46 |
| 500+ | ~0.47 |

The optimal value is **25 leaf nodes**, which was selected for the final model.

Random Forest variations (different estimators, criteria, depths) consistently produced MAE values around **0.43**, showing stability but no meaningful improvement over the optimized Decision Tree.

## 4) Feature Importance Analysis

Using Mutual Information scores, the following features were identified as the most influential in predicting sleep quality:

1. **Physical Activity Level**
2. **Stress Level**
3. **Heart Rate**
4. **Sleep Duration**
5. **Daily Steps**

These findings align with established sleep science studies, confirming the relevance of lifestyle factors in modeling sleep outcomes.

## 5) Final Selected Model

Based on the evaluation metrics, hyperparameter tuning, and performance stability:

- **Final Model:** Decision Tree Regressor (max_leaf_nodes = 25)

- **Final MAE:** ~0.42
- **Final RMSE:** ~0.60
- **R² Score:** ~0.65

# Conclusion

- The Decision Tree Regressor (max_leaf_nodes = 25) achieved optimal performance with an MAE of 0.42, while Random Forest models provided consistent results (MAE ~0.43), demonstrating the effectiveness of tree-based methods for sleep quality prediction.

- Key predictors—Physical Activity Level, Stress Level, Sleep Duration, Heart Rate, and Daily Steps—align with established sleep medicine research, validating our feature engineering approach and mutual information analysis findings.

- The hybrid prediction system, combining machine learning with domain knowledge adjustments for 10 individual factors (including age, activity level, stress, and physiological metrics), enables personalized sleep recommendations tailored to specific health profiles and lifestyle patterns.

- Future work includes integrating continuous sleep monitoring data, expanding the model to account for seasonal and environmental factors, incorporating cognitive performance metrics, and developing smartphone integration for real-time sleep optimization recommendation