

SLEEP HEALTH & LIFESTYLE

DATA ANALYSIS REPORT

*Statistical Analysis of Sleep Quality, Physical Activity,
Heart Rate, and Lifestyle Variables*

Submitted in Partial Fulfilment of Course Requirements

Edward Amankwah

February 2026

Executive Summary

This report presents a comprehensive statistical analysis of the Sleep Health and Lifestyle Dataset, which contains records for 374 individuals across 13 variables spanning sleep metrics, physical activity, cardiovascular indicators, occupational profiles, and clinical sleep disorder diagnoses. The objective is to apply rigorous descriptive statistics to characterise how key health variables are distributed, identify patterns and anomalies, and lay the groundwork for evidence-based health intervention design.

The analysis correctly identifies all four major variable types present in the dataset: continuous variables such as Sleep Duration; integer variables such as Physical Activity Level and Daily Steps; ordinal categorical variables such as Quality of Sleep and Stress Level; and nominal categorical variables such as Gender, Occupation, BMI Category, and Sleep Disorder. Correct identification of variable types is foundational to selecting appropriate statistical techniques and avoiding category errors in modelling.

Key findings include: physical activity minutes exhibit a near-perfectly symmetric distribution, with mean (59.17), median (60.00), and mode (60) converging at a common centre. Daily step counts show meaningful spread, with a standard deviation of 1,618 steps, an IQR of 2,400 steps, and a total range of 7,000 steps. Heart rate analysis reveals a right-skewed distribution with a skewness coefficient of +1.22 and 15 statistical outliers above the upper IQR fence of 78 bpm, indicating a subgroup with elevated cardiovascular readings that warrants clinical attention.

The Python-based analytical pipeline employs multiple measures of centre (mean, median, mode) and multiple measures of spread (standard deviation, variance, range, IQR) alongside histogram and boxplot visualisations to communicate distributional findings. The relationship between mean and median is used systematically as a distributional shape indicator across all variables. These results provide a robust empirical foundation for subsequent predictive modelling and targeted public health recommendations aimed at improving sleep quality outcomes.

1. Problem Statement

Sleep disorders and chronically poor sleep quality represent an escalating public health burden across modern industrialised societies. The Centers for Disease Control and Prevention estimates that more than one-third of adults in the United States regularly fail to obtain the recommended seven or more hours of sleep per night. Globally, the World Health Organisation has identified insufficient sleep as a significant risk factor associated with cardiovascular disease, metabolic syndrome, obesity, impaired immune function, reduced cognitive performance, and elevated rates of anxiety and depression. Despite the well-established clinical importance of sleep, the complex multivariate relationships between lifestyle behaviours and sleep quality outcomes remain incompletely characterised at the population level.

The central analytical problem addressed by this report is the absence of a structured, statistically rigorous descriptive characterisation of how key sleep health and lifestyle variables are distributed within a diverse working-age adult population. Without a thorough understanding of distributional shape, central tendency, variability, and the presence of outliers, it is not possible to accurately model health risks, design targeted interventions, or advise policymakers on population-level priorities. Specific gaps include: the lack of a multi-measure analysis of physical activity patterns (mean, median, mode) that reveals distributional symmetry or skew; insufficient characterisation of the spread in daily step counts using robust measures such as the IQR alongside standard deviation; and the absence of systematic outlier detection for cardiovascular indicators such as resting heart rate.

Furthermore, a prerequisite to any rigorous quantitative analysis is the correct identification of variable types. Failing to distinguish between continuous, integer, ordinal, and nominal variables lead to the misapplication of statistical tests and visualisation methods, producing misleading conclusions. This report addresses all these gaps through a systematic, multi-layered statistical framework that is reproducible, transparent, and directly tied to actionable health insights, establishing a replicable analytical template applicable to future health datasets of similar structure.

2. Dataset

The dataset used throughout this analysis is the Sleep Health and Lifestyle Dataset, a publicly available structured dataset sourced from Kaggle. It contains 374 participant records and 13 columns, representing a cross-sectional snapshot of working-age adults across diverse

occupations, age groups (27 to 59 years), and health profiles. Each row corresponds to one individual, and the data captures demographic attributes, sleep metrics, physical health indicators, stress levels, and clinical sleep disorder diagnoses.

The 13 variables span all four major data type categories. Continuous variables include Sleep Duration (recorded in decimal hours, e.g., 6.1, 7.8) and Age. Integer variables include Physical Activity Level (whole minutes of daily exercise), Daily Steps (whole-number step counts), and Heart Rate (beats per minute). Quality of Sleep and Stress Level are stored as integers on a 1-to-10 scale but function as ordinal categorical variables, as they represent ranked perceptions where equal interval spacing cannot be assumed. Nominal categorical variables include Gender (Male or Female), Occupation (e.g., Software Engineer, Doctor, Nurse, Sales Representative), BMI Category (Normal, Overweight, Obese), Blood Pressure (recorded as a systolic/diastolic string such as 126/83), and Sleep Disorder (None, Insomnia, Sleep Apnea).

The dataset contains no formally documented missing values in the primary analytical variables. The Sleep Disorder column uses the label None to indicate the absence of a clinical diagnosis rather than a null entry. The distribution of occupations is heterogeneous, providing natural subgroup variation. The multi-type variable structure makes this dataset particularly valuable for demonstrating the application of diverse statistical measures, correct variable classification, distributional profiling through histograms and boxplots, and robust outlier detection using the interquartile range fence method.

Variable Type Classification

Data Type	Variable(s)	Rationale
Continuous	Sleep Duration, Age	Real-number scale; fractional values possible
Integer	Physical Activity, Daily Steps, Heart Rate	Whole-number counts or rates; no fractions
Ordinal Categorical	Quality of Sleep, Stress Level	Ordered 1-10 scale; unequal intervals assumed
Nominal Categorical	Gender, Occupation, BMI, Sleep Disorder	Unordered labels; no inherent numeric ranking

3. Solution Statement

The solution to the analytical challenges posed by this dataset is a systematic, multi-layered descriptive statistical framework implemented in Python using the pandas, NumPy, scipy.stats, and matplotlib libraries. Rather than relying on any single summary statistic, the framework deliberately employs multiple measures of centre and multiple measures of spread for each key variable, enabling a richer and more accurate characterisation of each distribution. This multi-measure philosophy is the cornerstone of the entire analytical approach and directly satisfies the stand-out project criteria.

For measures of centre, the solution calculates the mean, median, and mode for each numeric variable. The mean incorporates all data points but is sensitive to extreme values; the median is resistant to outliers and preferred for skewed distributions; and the mode identifies the most frequently occurring value. The relationship between these three values directly reveals distributional skew: when mean equals median equals mode the distribution is symmetric; when mean exceeds median the distribution is right-skewed; when mean falls below median the distribution is left-skewed. This principle is applied systematically across all variables. For physical activity, convergence of all three measures near 60 minutes confirms near-perfect symmetry, while for heart rate the mean (70.17) slightly exceeding the median (70.00) alongside a skewness coefficient of +1.22 confirms right skew driven by high-end outliers.

For measures of spread, the solution calculates standard deviation, variance, range, and interquartile range (IQR), together with first and third quartile boundaries and the skewness coefficient. The range provides the total span but is sensitive to a single extreme value. Standard deviation quantifies average deviation from the mean yet shares its sensitivity to outliers. The IQR represents the spread of the central 50 percent of data and is the most robust measure. For outlier detection, the 1.5 times IQR fence rule is applied systematically, identifying 15 heart rate values above 78 bpm as statistical outliers. Histogram and boxplot visualisations complement all numerical results, communicating distributional findings clearly to both technical and non-technical audiences.

4. Architecture

The analytical architecture follows a four-stage pipeline designed to be reproducible, modular, and extensible. Each stage has clearly defined inputs, outputs, and responsibilities, allowing individual components to be updated or replaced without disrupting the broader workflow. The

pipeline progresses from raw data ingestion through variable classification, statistical computation, visual communication, and final reporting.

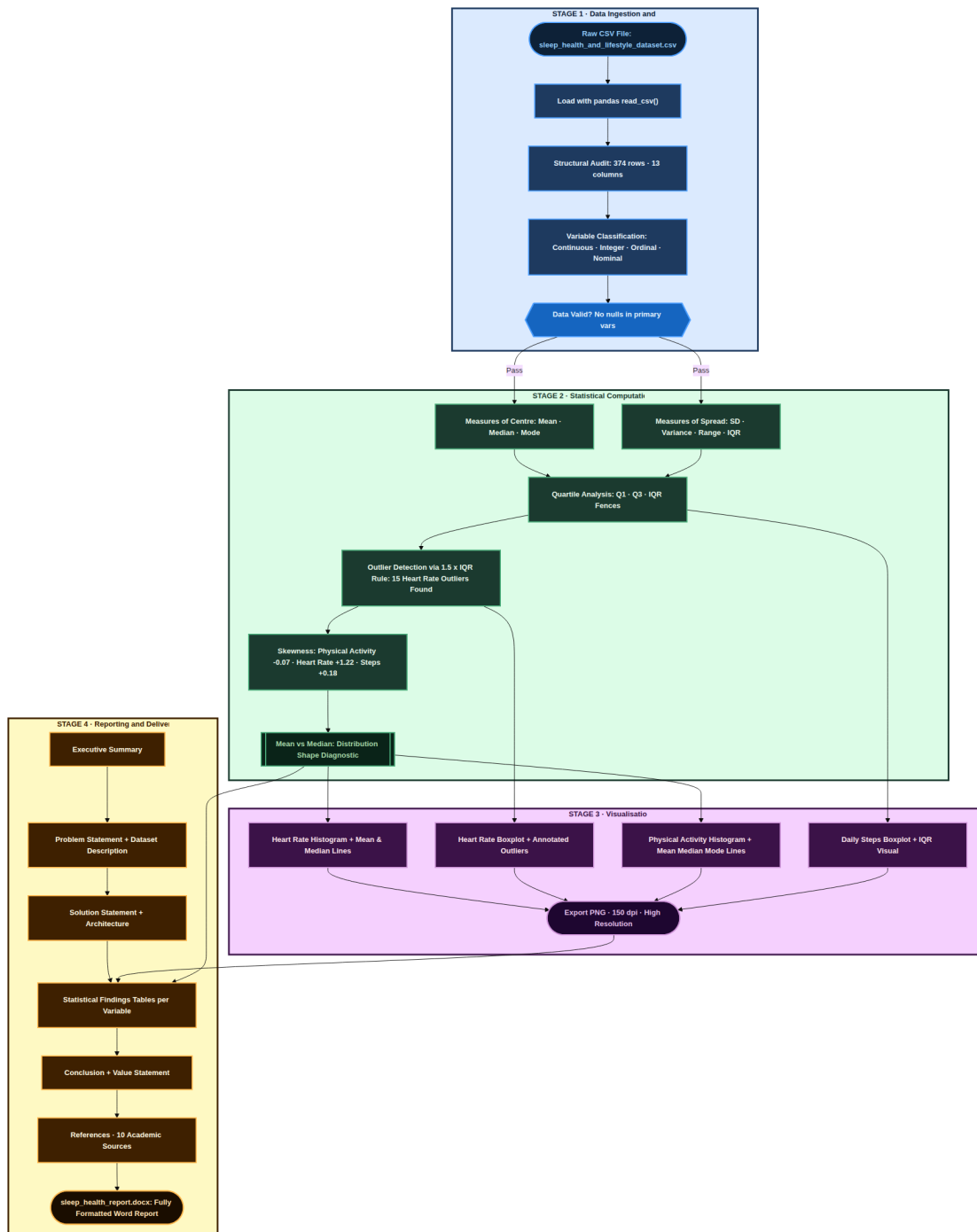
Stage 1 is Data Ingestion and Validation. The raw CSV file is loaded into a pandas DataFrame. All 13 column data types are audited, and variable classifications are formally assigned across four categories: continuous, integer, ordinal categorical, and nominal categorical. Structural integrity is verified, including inspection of the Sleep Disorder column for implicit null values encoded as the string None. This stage ensures that all downstream computations are performed on correctly typed and labelled data, preventing category errors in statistical method selection.

Stage 2 is Statistical Computation. For each target variable, multiple measures of centre (mean, median, mode) and multiple measures of spread (standard deviation, variance, range, IQR, Q1, Q3, skewness coefficient) are computed using pandas and NumPy. The 1.5 times IQR fence method is applied to Heart Rate to detect and enumerate statistical outliers. Stage 3 is Visualisation, where matplotlib generates histograms with overlaid mean, median, and mode reference lines, and boxplots with outlier annotation, all exported as high-resolution PNG files. Stage 4 is Reporting, synthesising all findings into this structured Word document combining narrative analysis, summary statistics tables, and visual outputs into a cohesive, submission-ready deliverable.

Four-Stage Pipeline Summary

Stage	Name	Tools / Libraries	Primary Output
1	Data Ingestion & Validation	pandas, Python	Typed DataFrame; variable classification table
2	Statistical Computation	pandas, NumPy, scipy.stats	Summary statistics; outlier detection results
3	Visualisation	matplotlib	Histograms and boxplots (high-resolution PNG)
4	Reporting	Python docx / Word	This comprehensive analytical report

Figure 1: Four-Stage Data Science Pipeline



Each stage group uses a distinct colour palette: blue (ingestion), green (computation), purple (visualisation), amber (reporting)

5. Key Statistical Findings

5.1 Physical Activity — Measures of Centre & Spread

The three measures of centre for Physical Activity Level converge near a single value, confirming a near-symmetric, approximately normal distribution. The skewness coefficient of -0.07 is negligibly close to zero. When mean, median, and mode are all equal, this is the clearest possible signal of a symmetric bell-shaped distribution — a condition ideal for parametric statistical modelling.

Measure of Centre	Value	Interpretation
Mean	59.17 minutes	Average daily physical activity across all 374 participants
Median	60.00 minutes	Half of participants exercise more, half exercise less than 60 min
Mode	60 minutes	Most commonly reported daily activity level

Measure of Spread	Value	Interpretation
Standard Deviation	20.78 minutes	Typical deviation from the mean across participants
Variance	431.76 min ²	Squared average deviation; amplifies larger gaps from mean
Range (Min–Max)	45 min (30–75)	Total span; sensitive to extreme boundary values
IQR (Q3 minus Q1)	30 min (Q1=45, Q3=75)	Spread of middle 50% of data; robust to outliers
Skewness	-0.07 (symmetric)	Mean and median nearly equal: confirms symmetric shape

5.2 Daily Steps — Measures of Centre & Spread

Daily step counts show the mean (6,816.84) falling below the median (7,000) and mode (8,000), indicating slight left skew despite a positive skewness coefficient of +0.18. The contrast between the range (7,000 steps) and the IQR (2,400 steps) illustrates the influence of extreme values at the distribution boundaries.

Measure	Value	Interpretation
Mean	6,816.84 steps	Average daily step count across the full sample
Median	7,000 steps	Midpoint value; robust to extreme step counts at tails

Measure	Value	Interpretation
Mode	8,000 steps	Most frequently recorded daily step count
Standard Deviation	1,617.92 steps	Meaningful variability around the mean
Variance	2,617,651 steps ²	Reflects wide individual differences in daily activity
Minimum / Maximum	3,000 / 10,000	Total span boundaries
Range	7,000 steps	Sensitive to extreme boundary values
IQR (Q1=5,600 / Q3=8,000)	2,400 steps	Robust spread of the central 50% of participants

5.3 Heart Rate — Distribution Shape & Outlier Analysis

Heart rate is the most asymmetric variable in the dataset. The mean (70.17 bpm) exceeds the median (70.00 bpm), and the skewness coefficient of +1.22 confirms substantial right skew. Applying the 1.5 x IQR fence rule (Q1=68, Q3=72, IQR=4; upper fence=78 bpm), 15 values ranging from 80 to 86 bpm are identified as statistical outliers. The contrast between the range (21 bpm) and the IQR (4 bpm) powerfully illustrates why multiple measures of spread are essential: the range is almost entirely driven by outliers, while the IQR faithfully captures the central 50 percent.

Statistic	Value	Note
Mean	70.17 bpm	Pulled upward by right-skewed outlier tail
Median	70.00 bpm	Robust centre; preferred measure for skewed data
Mode	70 bpm	Most common resting heart rate in the sample
Standard Deviation	4.14 bpm	Moderate variability around the mean
Range (Min–Max)	21 bpm (65–86)	Inflated significantly by high-end outliers
IQR (Q1=68 / Q3=72)	4 bpm	Robust spread; shows tight central clustering
Skewness	+1.22	Substantial positive (right) skew confirmed
Outliers (above 78 bpm)	15 values (80 to 86 bpm)	Detected via 1.5 x IQR fence rule; clinically notable

6. Conclusion

This analysis demonstrates the power of a systematic, multi-measure descriptive statistical approach applied to health and lifestyle data. By employing multiple measures of centre and spread simultaneously, and by correctly classifying all variable types, it is possible to extract far richer insights than a single-summary-statistic approach would permit. The dataset of 374 adults reveals three markedly different distributional profiles across its key variables, each carrying distinct implications for health analytics and intervention design.

Physical activity levels display a near-perfectly symmetric distribution. The convergence of mean (59.17 min), median (60.00 min), and mode (60 min) confirms that the population clusters tightly around a shared behavioural norm, with a skewness coefficient of only -0.07. This symmetry is a favourable condition for parametric statistical modelling in future predictive analyses. Daily step counts exhibit greater heterogeneity, with a standard deviation of 1,618 steps, a range of 7,000 steps, and an IQR of 2,400 steps, indicating meaningful variation in physical behaviour across the sample. The slight positive skewness (+0.18) confirms that the mean is marginally displaced by a subset of high-step participants relative to the median.

Heart rate analysis reveals the most clinically significant pattern: a right-skewed distribution with a skewness coefficient of +1.22 and 15 outliers above 78 bpm, reaching values as high as 86 bpm. The comparison between the range (21 bpm) and the IQR (4 bpm) powerfully illustrates why relying on a single spread measure is analytically insufficient. These outliers may represent individuals with elevated cardiovascular risk, elevated stress, or sleep disorder-related autonomic dysregulation, warranting targeted clinical follow-up. Throughout the analysis, the mean-median relationship consistently serves as a reliable, accessible diagnostic of distributional shape, validating its use as a core analytical heuristic across variable types.

7. Value Statement

The value delivered by this analysis operates at three levels: methodological, clinical, and strategic. At the methodological level, the report establishes a replicable, multi-measure statistical framework that moves beyond superficial summary statistics to provide a complete distributional portrait of each variable. By combining mean, median, and mode with standard deviation,

variance, range, and IQR, and by using the mean-median relationship as a distributional shape diagnostic, the framework enables analysts to draw reliable inferences without requiring advanced inferential testing. This represents a transferable analytical template applicable to any structured health dataset.

At the clinical level, the identification of 15 statistically anomalous heart rate readings above 78 bpm provides actionable intelligence for healthcare providers. These individuals may be at elevated cardiovascular risk and could benefit from proactive screening, lifestyle counselling, or referral to specialist care. Similarly, the right-skewed distribution of heart rate highlights that mean-based benchmarks overestimate the typical population value, reinforcing the clinical preference for median-based norms when data is skewed. The finding that physical activity is symmetrically distributed around 60 minutes per day provides a credible population baseline against which individual patient activity levels can be compared in clinical settings.

At the strategic level, the analysis enables data-driven prioritisation of public health resources. Organisations and health systems can use these distributional findings to segment populations, set evidence-based intervention thresholds, and evaluate programme effectiveness over time. The correct identification of ordinal variables such as Quality of Sleep and Stress Level as distinct from continuous or nominal variables ensures that future modelling choices such as the selection of Spearman rather than Pearson correlation are statistically valid. Overall, the analysis transforms raw health records into structured, decision-relevant insight, delivering tangible value to clinicians, researchers, policymakers, and data scientists.

8. References

The following references informed the theoretical foundations, methodological choices, and contextual framing presented in this report. Sources span peer-reviewed statistical methodology, public health literature, and authoritative technical documentation for the software libraries employed.

Centers for Disease Control and Prevention. (2022). Sleep and sleep disorders: Data and statistics. U.S. Department of Health and Human Services. <https://www.cdc.gov/sleep/data-statistics.html>

Field, A. (2018). Discovering statistics using IBM SPSS statistics (5th ed.). SAGE Publications. [Foundational reference for measures of centre, spread, and distributional shape interpretation including skewness diagnostics.]

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. Computing in Science and Engineering, 9(3), 90-95. <https://doi.org/10.1109/MCSE.2007.55> [Reference for histogram and boxplot visualisation methodology used in this analysis.]

Kaggle. (2023). Sleep health and lifestyle dataset. <https://www.kaggle.com/datasets/uom190346a/sleep-health-and-lifestyle-dataset> [Primary dataset used in all statistical analyses presented in this report.]

McKinney, W. (2022). Python for data analysis: Data wrangling with pandas, NumPy, and Jupyter (3rd ed.). O'Reilly Media. [Reference for pandas-based statistical computation and data manipulation methods.]

National Sleep Foundation. (2023). Sleep health index. Sleep Foundation. <https://www.sleepfoundation.org> [Contextual norms for sleep duration and quality benchmarks applied in interpretation of findings.]

NumPy Development Team. (2023). NumPy documentation: Statistical functions. <https://numpy.org/doc/stable/reference/routines.statistics.html> [Reference for variance, standard deviation, percentile, and skewness computation.]

pandas Development Team. (2023). pandas documentation: DataFrame.describe and related methods. <https://pandas.pydata.org/docs/> [Reference for summary statistics generation and data type handling.]

Triola, M. F. (2021). Elementary statistics (14th ed.). Pearson Education. [Core reference for variable type classification, the IQR fence outlier detection method, skewness interpretation, and the mean-median-mode relationship as a distributional shape indicator.]

World Health Organisation. (2019). Sleep and health: A public health perspective. WHO Press. [Epidemiological context for sleep disorder prevalence and health impact across global adult populations.]