The dataset I have chosen for this project has 3,000 rows and 20 columns, which contain both numerical and categorical data related to stress detection. Here's a breakdown of areas that may need cleaning and preprocessing based on my quick review of the data:

**1. Identify Data Issues**

- **Incorrect Values:**
  - Verify ranges and thresholds for numerical columns like PSS_score, sleep_duration, and screen_on_time.
  - Cross-check whether personality traits (Openness, Conscientiousness, etc.) and scores fall within the expected ranges.
- **NULL Values:**
  - Confirm there aren't any missing values as the data summary indicates, but recheck for anomalies such as placeholders (-1 or 9999).
- **Weird Characters:**
  - Ensure that all columns have consistent formatting, especially wake_time and sleep_time, which would most likely need a time conversion.

**2. Transforming the Data**

- **Feature Engineering:**
  - Obtain additional features such as Sleep Efficiency = sleep_duration / (wake_time - sleep_time).
- **Normalization:**
  - Fix columns like skin_conductance and accelerometer for compatibility in possible machine-learning models.

- **Encoding:**

  - Convert categorical columns (if any in extended data) to numbered formats.

- **Datetime Handling:**

  - Analyze any time-related fields into usable datetime formats if needed.

## 3. Data Subsetting

- Evaluate if all 20 columns are necessary for analysis or if any redundant features can be dropped.

## 4. Future-Proofing

- Clean and structure data for compatibility with Pandas, Visual Studio Code, and machine learning frameworks.

- Prepare data for statistical methods or modeling (linear regression, classification models).

## 5. Data Integrity Checks

- Validate that calculated metrics like mobility_distance and mobility_radius align with logical spatial constraints.