# Behavioural Patterns In Fitbit Users

December 2023

# Contents

# 1  Introduction

In the contemporary age of unparalleled technological advancements, the widespread adoption of health-tracking devices, such as Fitbits, has become a common practice. These devices generate extensive datasets, providing a unique opportunity to understand individual health patterns and broader societal trends. This study aims to comprehensively analyze Fitbit data, focusing on uncovering intricate patterns in users' behaviors, particularly in the context of step count, its influence on calorie expenditure, user grouping based on activity patterns, and a comparison with World Health Organization (WHO) activity level recommendations.

## 1.1  Background Information

Despite initial skepticism, recent research, exemplified by studies like Dubey Dubey (2019) and Henriksen Henriksen (2021), has demonstrated the positive impact of fitness trackers on users' health. Addressing concerns of demographic bias, Henriksen's work emphasizes the accessibility of fitness trackers to diverse populations. However, questions about data accuracy persist, as highlighted by Sieber's study Sieber (2023), urging a closer examination of the reliability of recorded data.

## 1.2  Research Objectives and Motivation

Building on existing credibility, this study aims to investigate Fitbit data, analyze behavioral trends. The research delves into group-level activity patterns, calorie expenditure predictors, and user clustering based on activity levels. Additionally, it evaluates individual adherence to WHO-recommended activity levels.

Based on the findings of the analysis, fitness tracker companies and better fine tune their personalized fitness models. Understanding personal and societal behaviour patterns allows companies to specifically find times of lower activity and target them in order to increase fitness. Other than just companies, governments or student organizations can also use this data in order to organize events that will bring in more people (e.g organize events on days where people are more active hence more likely to join events) or target times where people are not as active to help bring overall fitness and sense of community up (e.g more activities during winter as most people tend to spend more time at home).

## 1.3  Methodology

To achieve these objectives, rigorous data cleaning processes were implemented to ensure dataset integrity. Visual tools, including box plots, line plots, count plots, correlation heat maps, PCA, and K-means clustering, were employed for a nuanced understanding of the complex interplay between sustained physical activity and calories burnt over a 30-day period.

### 1.4 Key Findings

The study revealed compelling group-level trends, including heightened activity levels at the beginning of the week. Clustering analysis identified four distinct user groups with varying activity preferences. Comparison with WHO recommendations indicated that a portion of users did not meet daily requirements on certain days.

## 2 Problem Formulation

The primary focus of this project is a meticulous examination and analysis of Fitbit data, aiming to unravel intricate patterns that correlate an individual's calorie expenditure with their level of physical activity. At the heart of the investigation lies the exploration of the complex relationship between calories burned and various activities, with the goal of uncovering insights that significantly contribute to our understanding of health and fitness dynamics.

Navigating the complexities of human behavior and the myriad factors influencing fluctuations in calorie expenditure presents a central challenge in this endeavor. The project endeavors to address this challenge by harnessing the extensive dataset collected by Fitbit devices, allowing for a detailed examination of activity patterns over time. Key questions that guide this exploration include investigating whether specific activity levels or patterns exert a more pronounced influence on calorie expenditure and determining whether maintaining consistent activity throughout the week leads to more favorable outcomes compared to sporadic bursts of intense exercise. Through this approach, the project seeks to shed light on the nuanced interplay between physical activity and calorie burn, offering valuable insights into effective strategies for health and fitness management.

## 3 Dataset Description

The dataset under investigation originates from responses collected through a broad survey conducted on Amazon Mechanical Turk. This survey spanned a period of approximately two months, from March 12, 2016, to May 12, 2016. Notably, data was gathered for each individual over a span of thirty-one days. The uniqueness of this dataset lies in its source—thirty willing Fitbit users who generously provided their personal tracker data for analysis. Distinguished by its detailed granularity, the dataset includes minute-level records encompassing vital health metrics such as physical activity, heart rate, and sleep patterns.

### 3.1 Data Identification and Parsing

Identification of each contributor's data is facilitated by an export session ID (found in column A) or timestamp (located in column B), enabling versatile and

precise parsing of information. Each row of data corresponds to an individual's health metrics on a particular day.

### 3.2  Data Modification and Feature Engineering

Table 2 displays the original data, with certain features removed due to an abundance of zero values that did not contribute to the analysis. Simultaneously, new features like TotalMins, TotalHours, TotalActiveMins, TotalActiveHours, ActivityPreference, and Cluster were introduced.

### 3.3  Feature Definitions and Thresholds

TotalMins was calculated by adding up VeryActiveMinutes, FairlyActiveMinutes, LightlyActiveMinutes, and SedentaryMinutes. TotalHours is TotalMins divided by 60. TotalActiveMins was calculated by adding VeryActiveMinutes, FairlyActiveMinutes, and LightlyActiveMinutes. TotalActiveHours is TotalActiveMinutes divided by 60. Clustering was performed after the visual analysis, and which cluster each user belonged to was added to the dataset. ActivityPreference is categorical data produced by calculating the proportion of VeryActiveMins compared to their TotalActiveMins. Thresholds for ActivityPreference were defined as follows:

| Activity Preference | Threshold Values |
|---|---|
| Light | $x \leq 0.2$ |
| Moderate | $0.2 \leq x \leq 0.5$ |
| High | $x \geq 0.5$ |

Table 1: Threshold Values for Activity Preference

### 3.4  Dataset Evolution

An outstanding characteristic of this dataset is its inherent diversity, stemming from the utilization of various Fitbit tracker models and the distinct tracking behaviors and preferences demonstrated by each user. The inclusion of different Fitbit tracker models enriches the dataset, as these models may capture and report data in unique ways. This diversity presents an opportunity for future exploration into the impact of device variability on recorded metrics.

### 3.5  Dataset Tables

The final dataset can be seen in Table 3, illustrating the evolution from the original dataset (Table 1) to the modified dataset (Table 2).

| Variable | Data Type |
| --- | --- |
| ActivityDate | Date |
| TotalSteps | Integer |
| TotalDistance | Float |
| TrackerDistance | Float |
| LoggedActivitiesDistance | Interger |
| VeryActiveDistance | Float |
| ModeratelyActiveDistance | Float |
| LightActiveDistance | Float |
| SedentaryActiveDistance | Float |
| VeryActiveMinutes | Integer |
| FairlyActiveMinutes | Integer |
| LightlyActiveMinutes | Integer |
| SedentaryMinutes | Integer |
| Calories | Integer |

Table 2: Original Dataset

| Variable | Data Type |
| --- | --- |
| ActivityDate | Date |
| TotalSteps | Integer |
| TotalDistance | Float |
| TrackerDistance | Float |
| VeryActiveDistance | Float |
| ModeratelyActiveDistance | Float |
| LightActiveDistance | Float |
| VeryActiveMinutes | Integer |
| FairlyActiveMinutes | Integer |
| LightlyActiveMinutes | Integer |
| SedentaryMinutes | Integer |
| Calories | Integer |
| Weekday | Categorical |
| TotalMins | Integer |
| TotalHours | Integer |
| TotalActiveMinutes | Integer |
| TotalActiveHours | Integer |
| ActivityPreference | Categorical |
| Cluster | Integer |

Table 3: Modified Dataset

# 4 Methods

## 4.1 Libraries Used

The implementation of various data analysis and visualization tasks was facilitated by the utilization of the following Python libraries:

- **NumPy:** NumPy was employed for efficient handling of numerical operations and manipulation of arrays.

- **Pandas:** Pandas played a central role in data manipulation, offering powerful tools for data cleaning, exploration, and analysis through DataFrame structures.

- **Matplotlib:** Matplotlib was utilized for creating a diverse range of static visualizations, including line plots, scatter plots, and bar charts.

- **Seaborn:** Seaborn, built on top of Matplotlib, was used for generating aesthetically pleasing and informative statistical visualizations.

- **Scikit-Learn:** Scikit-Learn provided essential functionalities for data preprocessing, scaling, clustering (K-Means), and other machine learning tasks.

## 4.2 Data Loading and Cleaning

The study commenced with the loading of the primary dataset, "dailyActivity_merged.csv,". To ensure data integrity, rows with all recorded values of 0 were removed, mitigating potential distortions caused by incomplete or inaccurate entries. Along with all outliers past the 95th percentile had their values set to the value of the 95th percentile to avoid inaccuracies.

## 4.3 Feature Engineering

Feature engineering was a pivotal step in enhancing the analytical potential of the dataset, involving the creation of several key features to provide a more comprehensive understanding of daily activities.

### 4.3.1 Temporal Transformation

To facilitate temporal analysis, the 'ActivityDate' feature underwent a transformation into a datetime object. This conversion allowed for nuanced exploration of temporal patterns, enabling the identification of trends, cycles, and dependencies related to time.

### 4.3.2 Day-of-Week Categorization

A 'Weekday' column was introduced to categorize each record based on the corresponding day of the week. This addition served to capture potential variations in activity patterns based on weekdays, enabling the identification of trends related to workdays versus weekends and facilitating more granular analyses of weekly activity fluctuations.

### 4.3.3 Aggregated Metrics

Several aggregated metrics were computed to offer a holistic perspective on daily activity:

- **TotalMins:** The sum of VeryActiveMinutes, FairlyActiveMinutes, LightlyActiveMinutes, and SedentaryMinutes, providing a comprehensive measure of the total time spent in various activity levels.

- **TotalHours:** The total minutes converted into hours, offering a more user-friendly time unit for analysis and interpretation.

- **TotalActiveMins:** The sum of VeryActiveMinutes, FairlyActiveMinutes, and LightlyActiveMinutes, representing the total time spent in active pursuits.

- **TotalActiveHours:** The total active minutes converted into hours, providing a standardized unit for assessing active time.

These aggregated metrics allowed for a more nuanced evaluation of overall daily activity, considering both sedentary and active behaviors.

### 4.3.4 ActivityPreference Feature

The 'ActivityPreference' feature was derived to quantify the proportion of time spent in very active activities relative to total active minutes. This metric aimed to capture individual preferences for high-intensity activities within their overall active time, providing insights into the intensity distribution of activities. The 'ActivityPreference' feature facilitated the exploration of user tendencies toward specific activity levels and offered a finer-grained understanding of the composition of their active minutes. Refer to Table 1 for threshold values.

## 4.4 Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is a critical phase in data analysis, serving as a foundational step to understand the characteristics of the dataset. During EDA, various visualization techniques were employed to unravel patterns and trends in the data:

- **Box Plots, Histograms, and Visualizations:** Box plots, histograms, and other visualizations were utilized to provide a visual summary of the distribution and variability of key metrics. These graphical tools offer insights into central tendency, spread, and the presence of outliers.

- **Scaling Data using StandardScaler:** Subsequent analyses involved scaling the data using `StandardScaler`. This transformation standardizes features, ensuring a mean of 0 and a standard deviation of 1. This step facilitates comparisons between variables with different scales.

- **Visualizing Box Plots for Standardized Features:** After scaling, box plots for standardized features were generated. These visualizations enhance the understanding of feature distributions, especially after normalization, aiding in the identification of patterns and outliers.

- **Specific Insights from Visualizations:** Special attention was given to generating box plots for specific features, such as 'TotalSteps,' 'SedentaryMinutes,' and 'Calories.' These visualizations were analyzed across weekdays to uncover patterns in weekly activity. For example:

  - The 'TotalSteps' box plot could reveal variations in daily step counts throughout the week, highlighting days of increased or decreased activity.
  - 'SedentaryMinutes' distribution across weekdays might indicate patterns of sedentary behavior or periods of increased inactivity.
  - 'Calories' distribution insights across weekdays could provide information on variations in daily calorie expenditure.

### 4.5   Clustering

K-means clustering was applied to identify underlying patterns in the dataset. K-means clustering was used as it is one of the fastest and easiest ways to perform clustering and since just by looking at the data there were clear clusters a simple clustering models suffices. The parameters for the clustering were max_iter = 300, n_init = 10, and random_state = 42.The Elbow Method guided the selection of an optimal number of clusters which ended up being 4. The addition of a 'Cluster' column to the original dataset allowed for subgroup analyses based on cluster assignments.

### 4.6   Meets Daily Requirements

A binary classification was introduced to categorize individuals based on whether they met the daily recommendation of at least 20 minutes of total active minutes. This binary indicator, 'Meets_Daily_Requirements,' served as a fundamental metric for evaluating adherence to daily activity guidelines.

# 5 Results

The visual analysis of feature distributions indicates that, for the most part, the dataset exhibits characteristics resembling a normal distribution.
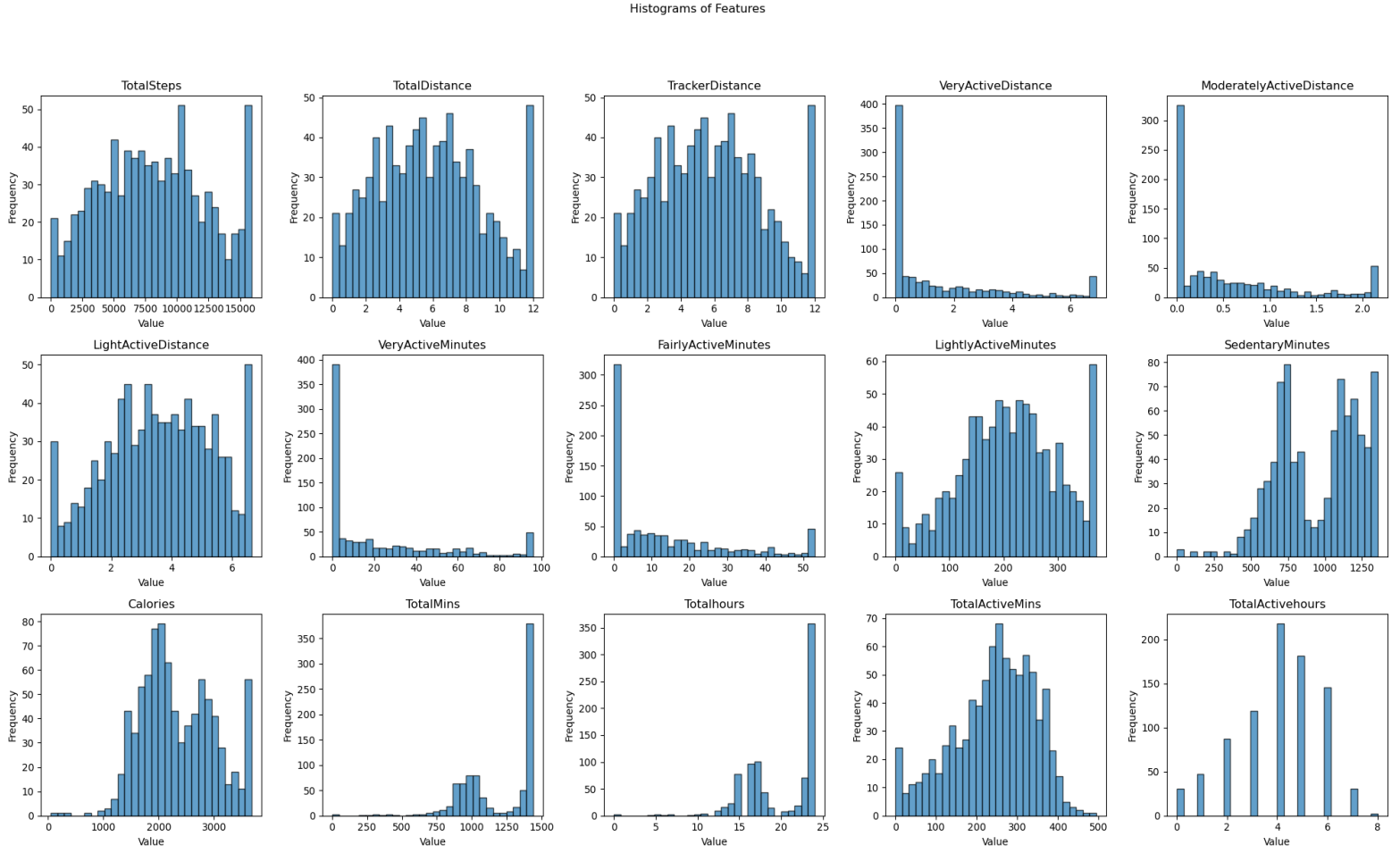


Figure 1: Feature Distributions Before Scaling

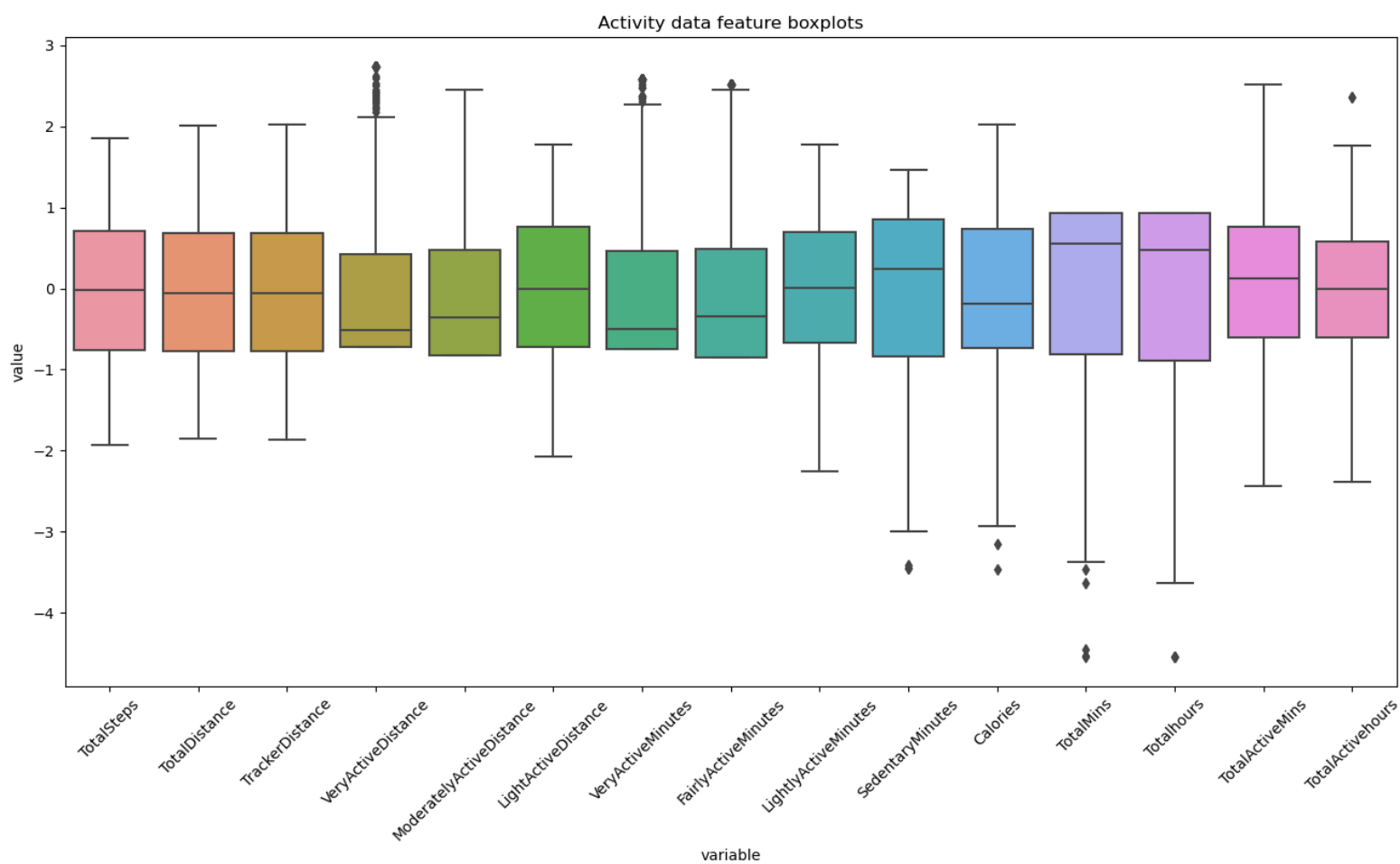After scaling, the distribution of features takes on a different appearance, as demonstrated below:

Figure 2: Feature Distributions After Scaling

Specifically observing the distributions of TotalSteps, SedentaryMinutes, and Calories per weekday:
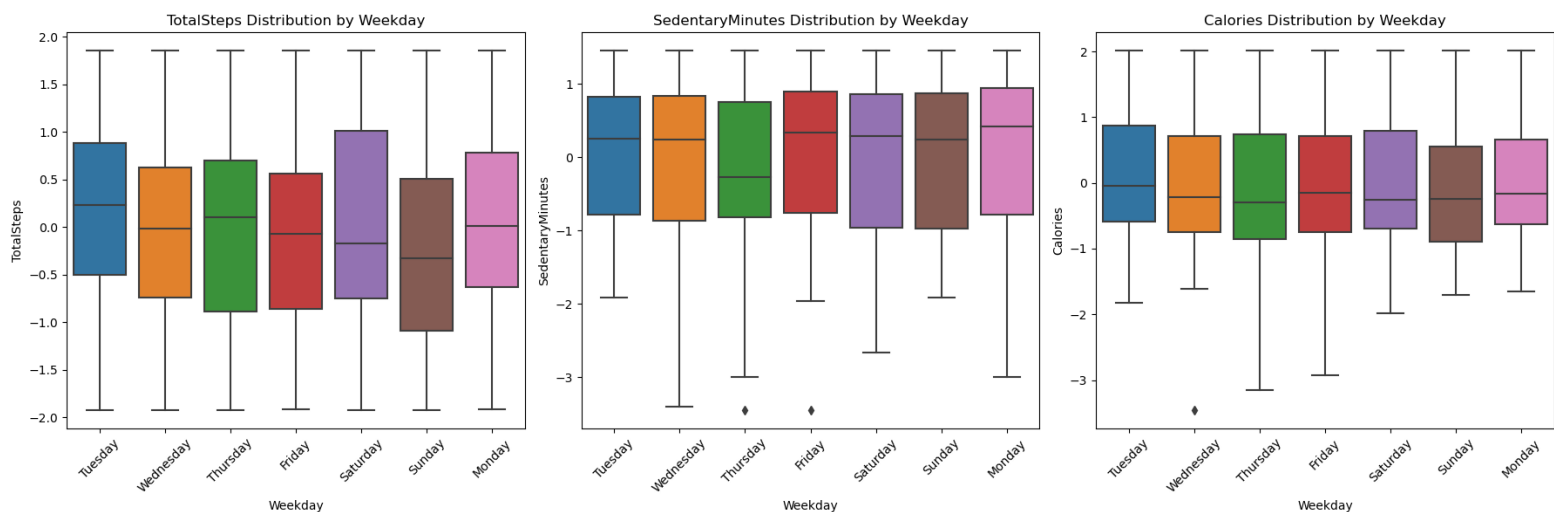
Figure 3: Weekday-wise Feature Distributions

Analysis of these plots reveals noticeable trends, notably, an increase in TotalSteps and Calories on Mondays, Tuesdays, and Saturdays, while experiencing a sharp decline on Wednesdays and especially on Sundays.

The correlation heatmap below showcases the key predictors for the amount of calories burnt:
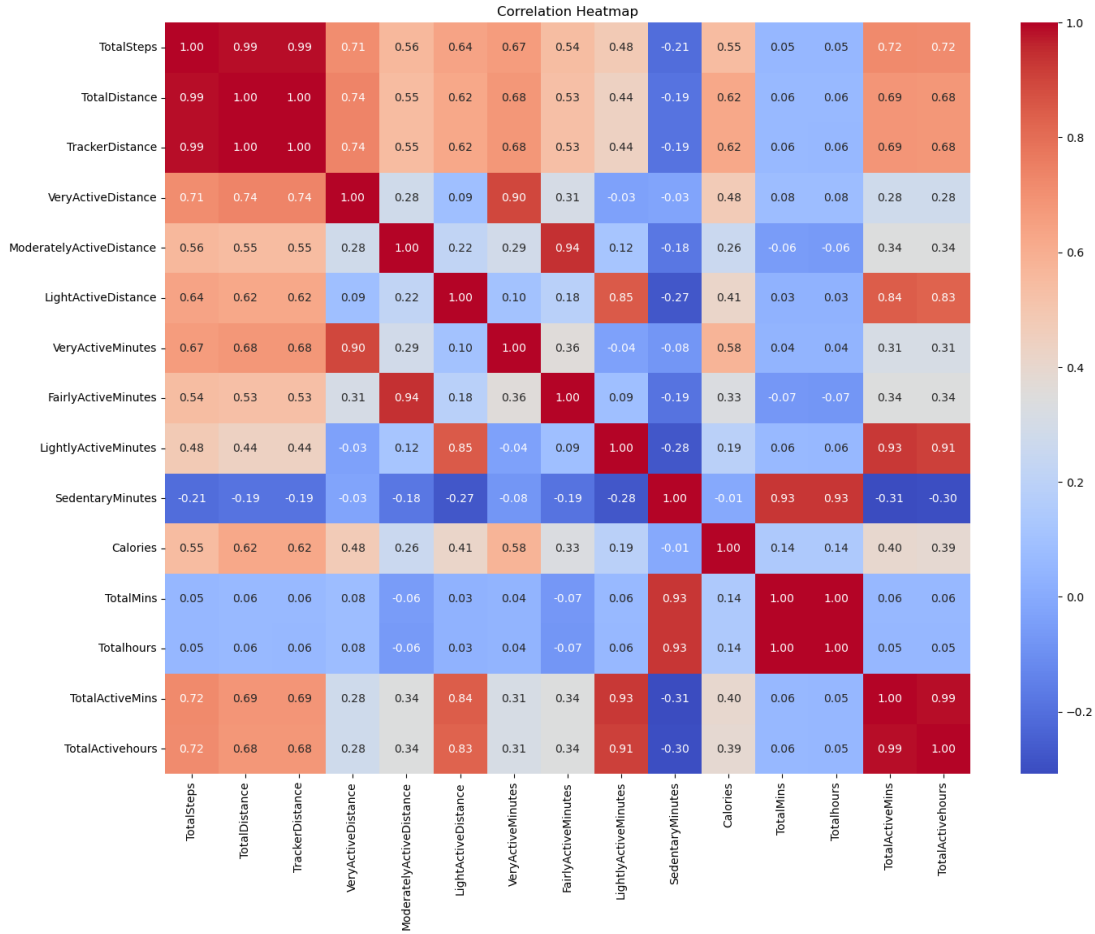
Figure 4: Correlation Heatmap of Predictors for Calories Burnt

From the heatmap, we can clearly see that the best predictors for the amount of calories burnt are TotalDistance and VeryActiveMinutes. These variables exhibit a strong positive correlation with calorie expenditure.

Further examination involves studying the behavioral patterns of a randomly selected individual from the group. The individual, id 4388161847, exhibits an above-average level of total activity, accompanied by lower VeryActiveMinutes. However, they demonstrate significantly higher LightlyActiveMinutes, resulting in a higher-than-average calorie expenditure.
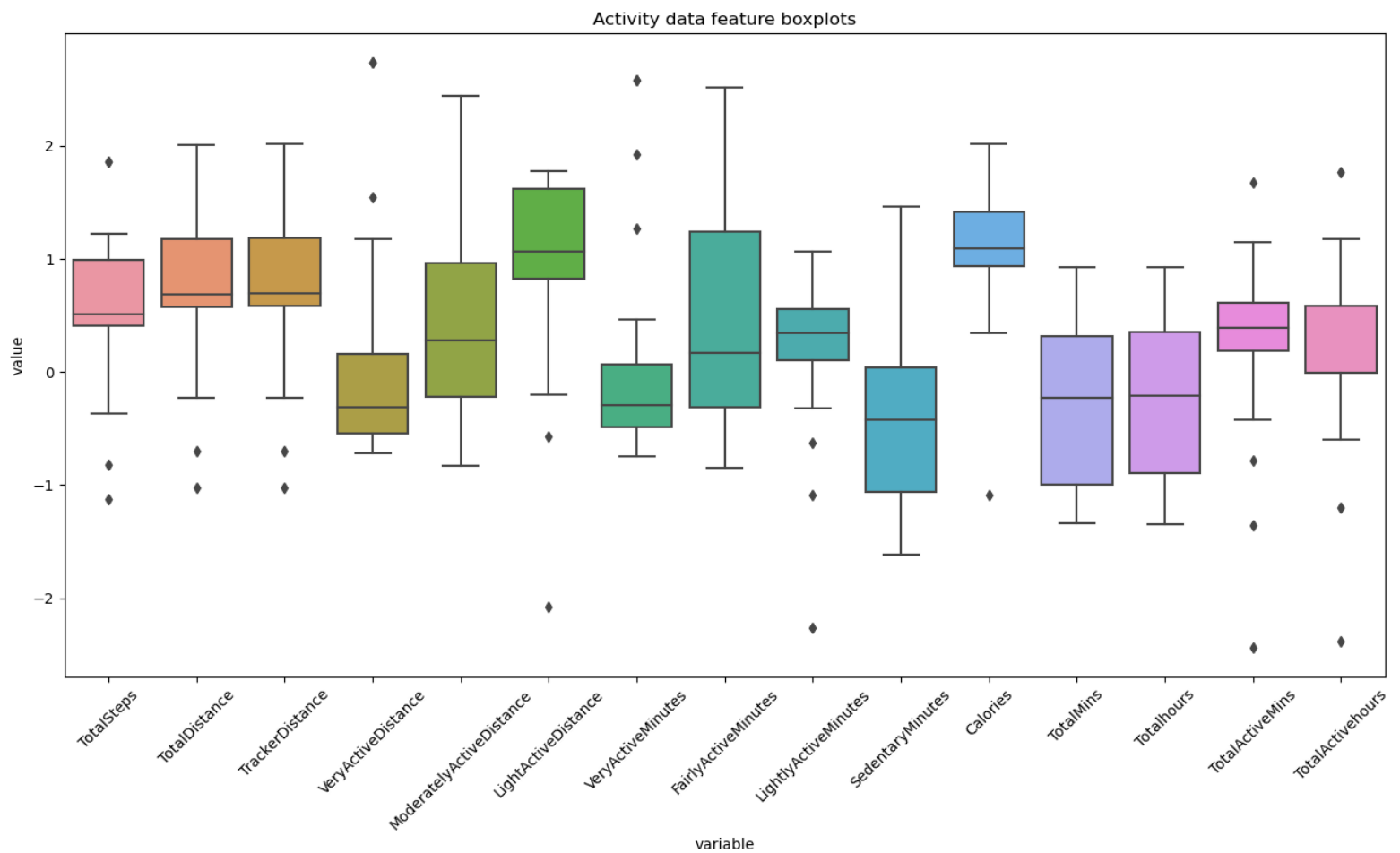
Figure 5: Behavioral Patterns of a Selected Individual

Finally, the clustering results, as depicted in Figure 6, provide valuable insights into the grouping of individuals based on their activity patterns. The clustering algorithm identified four distinct clusters within the dataset. Let's delve deeper into the characteristics of each cluster:
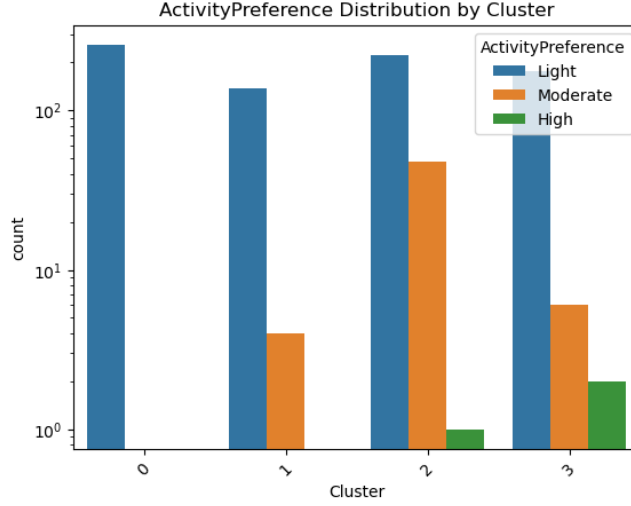
Figure 6: Cluster Distribution

| Cluster | Activity Characteristics |
|---------|--------------------------|
| 0 | Predominantly light activity. Members of this cluster exhibit lower levels of overall physical activity, focusing primarily on light activities. Their daily routines involve minimal moderate or vigorous activities. |
| 1 | Mostly light activity with a hint of moderate activity. Individuals in Cluster 1 engage in mostly light activity, similar to Cluster 0. However, this group includes a small proportion of moderate activity, indicating a slightly more varied activity profile compared to Cluster 0. |
| 2 | Mostly light and moderate activity with some high activity. Cluster 2 represents individuals who maintain a balance between light and moderate activity, with a notable inclusion of high-intensity activities. Members of this cluster exhibit a more diverse and active lifestyle compared to Clusters 0 and 1. |
| 3 | Mostly light activity with a significant proportion of high activity. Cluster 3 is characterized by individuals primarily involved in light activity but with a substantial proportion of high-intensity activity in their routines. This group showcases a unique combination, suggesting a more dynamic and physically active lifestyle compared to the other clusters. |

Table 4: Cluster Characteristics

Understanding these clusters allows for targeted insights into different activity patterns among participants. This information can be instrumental in tailoring health recommendations, interventions, or fitness programs to better suit the preferences and habits of individuals within each cluster. It also facilitates a nuanced understanding of how various factors contribute to different activity profiles, contributing to the broader discourse on personalized health and fitness strategies.

In accordance with the guidelines set by the World Health Organization (WHO), individuals are recommended to engage in at least 150 minutes of moderate-intensity aerobic physical activity throughout the week *Physical Activity* (n.d.). This translates to an average of 20 minutes of activity per day. Recognizing the importance of consistent physical activity in maintaining optimal health, the WHO's guidelines serve as a benchmark for promoting an active lifestyle.

Upon examining the dataset and evaluating the daily activity levels of the participants, it is noteworthy that the majority of individuals faced challenges in meeting the recommended 20 minutes of daily activity consistently. Specifically, the analysis reveals that 26 out of the 30 participants fell short of achieving this daily goal on at least one occasion. However, it is encouraging to observe that despite these occasional lapses, all participants successfully accumulated sufficient weekly activity to meet or exceed the WHO's overall weekly recommendation of 150 minutes

# 6  Conclusion and Discussion

## 6.1  Summary of Findings

The comprehensive analysis of Fitbit data has provided valuable insights into users' behaviors, physical activity patterns, and calorie expenditure. The study began with an exploration of feature distributions, revealing characteristics akin to a normal distribution. Scaling these features and analyzing specific metrics across weekdays unveiled trends, such as increased activity on certain days and declines on others.

The correlation heatmap highlighted TotalDistance and VeryActiveMinutes as crucial predictors for calorie expenditure. This finding reinforces the importance of both overall activity levels and high-intensity exercise in influencing energy burn. The examination of an individual's behavioral patterns further emphasized the complexity of responses to activity, highlighting the need for personalized health strategies.

Clustering analysis identified four distinct clusters with varying activity profiles. From predominantly light activity to a combination of light, moderate, and high-intensity activities, understanding these clusters facilitates tailored health recommendations. The introduction of 'ActivityPreference' as a metric quantifying the proportion of time spent in very active activities added granularity to the analysis, capturing individual intensity preferences.

## 6.2 Implications and Recommendations

Based on the findings of the analysis, the implications for fitness tracker companies are significant. The ability to fine-tune personalized fitness models by understanding individual and societal behavior patterns opens new avenues for enhancing the efficacy of fitness tracking devices. Companies in the fitness technology sector can leverage these insights to pinpoint times of lower activity among users and strategically target those periods to encourage increased physical activity. This targeted approach could include personalized notifications, challenges, or incentives to motivate users during times when they are less active.

In the context of personalized health and fitness management, the study emphasizes the importance of tailoring recommendations based on individual preferences and habits, as highlighted by the clustering analysis. Acknowledging and accommodating diverse activity patterns among users can significantly enhance the impact of health and fitness interventions. Fitness tracker companies, health professionals, and policymakers can collaboratively work towards developing personalized strategies that resonate with users, ultimately fostering long-term adherence to physical activity guidelines.

## 6.3 Limitations and Challenges

While the study has provided valuable insights into Fitbit users' behaviors and activity patterns, it is essential to recognize and address several limitations that may influence the interpretation and generalization of the findings. The dataset's origin from responses collected through Amazon Mechanical Turk introduces a potential source of sampling bias. Participants recruited through this platform may not be fully representative of the broader population, as Mechanical Turk users may possess distinct characteristics or motivations that differ from the general public. This sampling bias can limit the external validity of the study, and caution is warranted when generalizing the findings to populations beyond the sampled group.

The study's reliance on a sample size of thirty users, while providing valuable insights, may not capture the full diversity of populations using fitness trackers. A larger and more diverse sample would enhance the study's robustness and allow for more nuanced subgroup analyses. The findings, therefore, should be considered within the context of the study's specific sample and may not be universally applicable. The dataset's accuracy raises concerns, as highlighted by previous studies. While efforts were made to address this through rigorous data cleaning processes, inherent inaccuracies in self-reported fitness tracker data remain a challenge. Users may vary in their diligence in reporting, and discrepancies between self-reported and actual activity levels could impact the reliability of the results.

The binary classification of individuals based on whether they meet daily activity requirements, though providing a clear metric for analysis, oversimplifies the complex nature of human behavior. The WHO-recommended thresholds,

while widely accepted, may not universally apply to all individuals. Individuals have diverse lifestyles, health conditions, and preferences, making it challenging to categorize adherence to daily activity guidelines in a binary manner.

The study acknowledges this limitation and encourages a nuanced understanding of individual variability. The decision to remove all records with 0 values introduces the possibility of bias in the results. Individuals who predominantly engage in sedentary behaviors or have periods of inactivity may be underrepresented in the dataset due to the removal of these records. This exclusion may impact the generalizability of the findings, especially when assessing patterns related to sedentary behavior or intermittent activity.

Estimating calorie expenditure based on activity metrics involves inherent challenges. Individual variations in metabolism, body composition, and other factors may influence the accuracy of calorie expenditure calculations. The study acknowledges the complexities of calorie expenditure estimation and the need for additional factors to be considered in future analyses.

Addressing these limitations and considering them in the interpretation of the results is crucial for maintaining the integrity and applicability of the study's findings. Future research endeavors should aim to overcome these challenges, potentially through larger and more diverse datasets, improved data accuracy validation methods, and a more nuanced approach to activity classifications.

### 6.4   Future Directions

Future research endeavors should address these limitations and explore additional dimensions of individualized health strategies. This includes:

1. **Diverse Demographics:** Conduct studies with larger and more diverse populations to ensure broader applicability of findings and mitigate demographic biases.

2. **Longitudinal Analysis:** Extend the analysis over a more extended period to capture long-term trends and variations in activity patterns.

3. **Accuracy Assessment:** Implement more robust methods to assess and improve the accuracy of recorded data, considering the limitations identified in previous research.

4. **Dynamic Thresholds:** Explore dynamic thresholds for activity preferences, acknowledging that individual preferences may evolve over time.

5. **Device-Specific Analyses:** Conduct in-depth analyses of the impact of device variability on recorded metrics, considering factors such as sensor accuracy and tracking algorithms.

### 6.5   Conclusion

In conclusion, the study delves into the intricate relationship between physical activity, calorie expenditure, and individual preferences. The findings contribute to the ongoing discourse on personalized health and fitness management,

emphasizing the need for tailored approaches. Despite limitations, the study's insights lay the groundwork for future research endeavors that can refine our understanding of health and contribute to more effective interventions.

# References

Dubey, D. (2019). *An analytical study of use and effects of fitness tracker on humans.* Retrieved from `https://www.researchgate.net/publication/342449570_AN_ANALYTICAL_STUDY_OF_USE_AND_EFFECTS_OF_FITNESS_TRACKER_ON_HUMANS`

Henriksen, A. (2021). *Using fitness trackers and smartwatches to measure physical activity in research: Analysis of consumer wrist-worn wearables.* Retrieved from `https://www.academia.edu/52967921/Using_Fitness_Trackers_and_Smartwatches_to_Measure_Physical_Activity_in_Research_Analysis_of_Consumer_Wrist_Worn_Wearables`

*Physical activity.* (n.d.). `https://www.who.int/news-room/fact-sheets/detail/physical-activity`. (Accessed: Insert Date Here)

Sieber, C. (2023). Title of the article. *Frontiers in Digital Health*, *1*, 1-10. Retrieved from `https://www.frontiersin.org/journals/digital-health/articles/10.3389/fdgth.2023.1006932/full`