

# **FitLifeBot : Personalized Fitness Recommendation Web App**

By

Keerthi Nalabothula - 002851756

Pavan Kumar Reddy Sajjala - 002794873

## **Abstract**

The exponential growth of wearable technology and health tracking applications has revolutionized how individuals monitor their physical activities. Despite the extensive data provided by platforms like Google Fit, users often struggle to extract meaningful and personalized insights to inform their fitness routines. The "FitLifeBot: Personalized Fitness Recommendation System" project addresses this gap by leveraging advanced machine learning and natural language processing techniques.

FitLifeBot processes Google Fit data to generate tailored fitness recommendations, enhancing the utility of health tracking applications. The system analyzes user data, identifies patterns, and provides insights grounded in users' unique activity patterns and health metrics. The project integrates a user-friendly interface for data input and output, converting raw data into actionable insights. This empowers users to make informed decisions about their fitness routines, promoting healthier lifestyles and improved fitness outcomes.

The development and deployment of FitLifeBot involved creating a web-based platform that seamlessly integrates with Google Fit. The platform includes a machine learning model for data analysis and an LLM for generating personalized recommendations. This innovative approach ensures users receive relevant advice to optimize their fitness journeys, bridging the gap between data availability and practical application.

# Table of contents

Abstract.....	2
<b>CHAPTER 1</b>	
<b>INTRODUCTION.....</b>	<b>1</b>
Chapter 1: Introduction.....	2
1.1 Background.....	2
1.2 Motivation.....	3
1.3 Problem Statement.....	3
1.4 Scope & Objectives.....	4
<b>CHAPTER 2</b>	
<b>LITERATURE REVIEW.....</b>	<b>5</b>
Chapter 2: Literature Review.....	6
2.1 Introduction.....	6
2.2 Related Literature on Fitness Tracking and Recommendations.....	6
2.3 Fitness Recommendation Systems and Thier Limitations.....	7
2.4 Innovations in FitlifeBot.....	7
2.5 Additional Studies.....	7
2.6 Summary.....	8
<b>CHAPTER 3.....</b>	<b>9</b>
<b>METHODOLOGY.....</b>	<b>9</b>
Chapter 3: Methodology.....	10
3.1 Introduction.....	10
3.2 Ethical and User-Centric Data Collection.....	10
3.2.1 Platform Rationale & Limitations:.....	10
3.2.2 Website Integration:.....	11
3.3 Advanced Machine Learning for Activity Analysis.....	12
3.4 Model Training and Validation.....	13
3.4.1 Feature Selection:.....	13
3.4.2 Train-Test Split.....	13
3.4.3 Pipeline Creation.....	14
3.4.4 Cross-Validation.....	14
3.4.5 Model Evaluation.....	14
3.4.6 PDF Report Generation.....	15
3.5 LLM as a Personalized Fitness Coach.....	15
<b>CHAPTER 4.....</b>	<b>17</b>
<b>RESULTS.....</b>	<b>17</b>
Chapter 4: Results.....	18
Results: Analysis of Human Activity Recognition Datasets for FitLifeBot Project.....	18
Analysis of Google Fit Data for FitLifeBot Project.....	24
Google Fit Data Analysis and Ensemble Model.....	24
Monthly User Fitness Report.....	27
Summary.....	28
<b>CHAPTER 5.....</b>	<b>30</b>
<b>Discussions.....</b>	<b>30</b>
Chapter 5: Discussions.....	31
Discussion: Analysis of Google Fit Data for FitLifeBot Project.....	31
<b>CHAPTER 6.....</b>	<b>35</b>

**CONCLUSION.....35**

Chapter 6: Conclusion..... 36

    Implications..... 37

    Limitations..... 37

    Future Work..... 38

Conclusion..... 39

References..... 40

# **CHAPTER 1**

## **INTRODUCTION**

# Chapter 1: Introduction

## 1.1 Background

In recent years, the proliferation of wearable devices and health-tracking applications has dramatically transformed the landscape of personal fitness management. These technologies have allowed individuals to collect extensive data about their physical activities and overall health. Platforms like Google Fit have become ubiquitous, allowing users to monitor various health metrics, such as steps taken, calories burned, heart rate, and sleep patterns. Despite the vast amount of data these applications generate, many users need help to utilize this information effectively due to its complexity and the need for personalized analysis.

The primary challenge lies in users' ability to translate raw data into meaningful insights to guide their fitness routines and lifestyle choices. While Google Fit and similar platforms provide a wealth of quantitative data, they often need more analytical tools to interpret this data in a personalized and actionable manner. This disconnect between data availability and practical utility highlights a significant gap in the fitness technology ecosystem.

Recognizing this gap, there is a pressing need for systems that can manage and process large volumes of fitness data and provide personalized recommendations based on individual user's unique needs and patterns. Integrating advanced technologies like machine learning (ML) and AI can significantly enhance our processes. Natural language processing (NLP) holds promise in addressing this need. ML can analyze complex datasets, identify patterns, and generate insights tailored to the user's health and fitness profile. On the other hand, NLP can be used to understand and interpret user queries and feedback, making the system more interactive and user-friendly.

The development of such a system requires a comprehensive approach, beginning with the collection and preprocessing of data to ensure its accuracy and relevance. This is followed by applying sophisticated ML algorithms that detect trends and patterns within the data. Additionally, an intuitive user interface enables users to interact with the system seamlessly and access personalized recommendations easily. The objective is to establish a comprehensive platform that simplifies data interpretation, relieving users from the burden of complex data, and enhancing the user's ability to make informed decisions about their health and fitness, thereby providing a sense of relief.

## 1.2 Motivation

The popularity of health and fitness applications reflects a growing societal trend towards health consciousness and self-improvement. However, the vast amount of data these applications collect can be overwhelming for users who need more expertise to interpret it effectively. There is a clear need for a system that tracks physical activities and provides personalized, easy-to-understand recommendations to guide users toward their fitness goals.

FitLifeBot addresses this need using sophisticated machine-learning models to analyze activity data and generate tailored fitness advice. Unlike generic fitness tips, the recommendations provided by FitLifeBot are based on the user's unique activity patterns and health metrics, ensuring relevance and effectiveness. This personalized approach empowers users to make informed decisions about their well-being and motivates them to stay active and healthy.

## 1.3 Problem Statement

Despite the wealth of data provided by Google Fit, users face several challenges in improving their fitness routines:

- Difficulty in interpreting raw activity data to derive personalized insights
- The challenge lies in providing personalized fitness recommendations that adapt to individual users' unique health and activity data.
- Inability to integrate and analyze data from multiple sources to provide a comprehensive fitness plan

These challenges underscore the need for an intelligent system capable of processing and analyzing fitness data to offer meaningful, personalized recommendations. FitLifeBot aims to fill this gap by providing users a comprehensive fitness plan, customized fitness advice based on their activity data, helping them achieve their health and fitness goals more effectively, and giving them a sense of control over their fitness journey.

## 1.4 Scope & Objectives

The proposed system will simplify and automate various processes:

### *Data Collection and Preprocessing:*

- Collect fitness data from Google Fit
- Cleanse, transform, and structure the raw data to optimize it for analytical processing and modeling

### *Model Training:*

- Model Training: The system will be trained to analyze fitness data and predict user needs. Sophisticated machine learning models and data analysis techniques will be implemented to refine and personalize recommendations, enhancing their precision and applicability.

### *User Interface:*

- Develop a user-friendly interface to display personalized fitness recommendations
- Ensure that the interface is intuitive and accessible to users of all levels of technical expertise

### *Evaluation:*

- Test the system with actual user data to evaluate its effectiveness
- Utilize feedback to enhance and perfect the system iteratively, ensuring continuous improvement and refinement

The primary objective of this project is to develop FitLifeBot, a personalized fitness recommendation system that utilizes machine learning and natural language processing to analyze user activity data and offer customized fitness advice. The specific objectives are:

- To collect and preprocess fitness data from Google Fit and other sources
- To train machine learning models to recognize patterns in the data and predict user needs
- To develop an interface that provides users with personalized fitness recommendations based on their activity data
- To evaluate the effectiveness of the recommendations in improving user fitness levels



# **CHAPTER 2**

## **LITERATURE REVIEW**

## Chapter 2: Literature Review

### 2.1 Introduction

The purpose of this chapter is to review past efforts related to the proposed FitLifeBot system, highlighting their limitations and identifying gaps that necessitate our solution. Numerous studies have explored the fields of fitness tracking, machine learning in health, and personalized recommendation systems. Fitness tracking involves collecting data about an individual's physical activities using wearable devices or mobile applications. Machine learning in health leverages algorithms to analyze this data and provide actionable insights. The literature reviewed in this chapter emphasizes the following points: advancements and shortcomings in these areas, setting the stage for our innovative approach.

### 2.2 Related Literature on Fitness Tracking and Recommendations

A comprehensive fitness tracking system provides users detailed data about their physical activities, health metrics, and progress toward fitness goals. Early fitness tracking systems relied heavily on manual data entry and simple statistical analyses. However, wearable technology has revolutionized this field, enabling real-time data collection and more sophisticated analyses. For instance, systems like Fitbit and Apple HealthKit track various activities and health metrics but often fall short of providing personalized insights that are actionable and tailored to individual users' needs (Smith et al., 2018).

Machine learning has emerged as a powerful tool in health and fitness, offering the ability to analyze large datasets and uncover patterns that can inform personalized recommendations. A study by Zhang et al. (2020) demonstrated the potential of using machine learning algorithms to predict user activity levels and provide customized exercise plans. Despite these advancements, many existing systems need more integration with popular fitness tracking platforms like Google Fit and fail to provide a seamless user experience.

Integrating natural language processing (NLP) with fitness tracking systems is relatively new. NLP can enhance user interaction by enabling conversational interfaces and providing more intuitive recommendations. Research by Kim et al. (2019) explored the use of chatbots in health applications, finding that they significantly improved user engagement and adherence to fitness routines. However, there remains a need for systems that combine machine learning, NLP, and real-time data integration to offer truly personalized and actionable fitness recommendations.

## 2.3 Fitness Recommendation Systems and Their Limitations

While fitness recommendation systems have become more prevalent, they often need more personalization and adaptability. Traditional systems provide generic advice based on broad user categories rather than individual data. Hernandez et al. (2021) highlighted the importance of personalized recommendations in improving user outcomes. Their study found that users who received tailored advice based on their specific activity patterns and health metrics were more likely to adhere to fitness plans and achieve their goals.

Despite these findings, many commercial systems still rely on static rules and predefined routines. For example, the MyFitnessPal app offers diet and exercise suggestions, often based on general guidelines rather than the user's unique data profile. This method could result in user dissatisfaction and reduced effectiveness of the recommendations.

## 2.4 Innovations in FitLifeBot

FitLifeBot aims to address the limitations identified in the literature by integrating advanced machine learning algorithms with real-time data from Google Fit. Unlike existing systems, FitLifeBot leverages a large language model (LLM) to provide natural language recommendations, enhancing user interaction and engagement. By continuously analyzing user data, the system can adapt recommendations based on evolving activity patterns and health metrics, ensuring they remain relevant and practical.

Our approach builds on the work of previous studies but pushes the boundaries by offering a more holistic and user-centric solution. Integrating machine learning, NLP, and real-time data collection creates a dynamic system capable of delivering highly personalized and actionable fitness advice.

## 2.5 Additional Studies

Several additional studies provide insights into the development of advanced fitness tracking and recommendation systems. For instance, a study by Chen et al. (2019) on integrating IoT and machine learning technologies in wearable devices highlights the potential for real-time monitoring and dynamic recommendations. Their findings suggest real-time feedback significantly improves user engagement and fitness outcomes (Chen et al., 2019).

Furthermore, research by Li and Bai (2020) on adaptive fitness coaching systems emphasizes the

importance of continuous learning algorithms that evolve with user data. Their adaptive system demonstrated higher user satisfaction and effectiveness than static recommendation models (Li & Bai, 2020). This underscores the need for FitLifeBot's adaptive capabilities, ensuring it remains responsive to user needs and preferences over time.

## 2.6 Summary

The literature review shows an increasing focus on leveraging technology to promote health and fitness, mainly through fitness tracking and personalized recommendations. Wearable devices and mobile applications have revolutionized fitness tracking, enabling real-time data collection and analysis. Analyzing this data using machine learning algorithms has demonstrated potential in offering personalized recommendations and insights. Nevertheless, current solutions frequently require additional personalization, smooth integration with widely-used platforms, and the capacity to adjust to individual needs and preferences. FitLifeBot aims to address these limitations by integrating advanced machine learning models with real-time data from Google Fit and leveraging a large language model (LLM) to provide natural language recommendations. This approach builds on previous research but pushes the boundaries by offering a more holistic and user-centric solution that combines machine learning, NLP, and real-time data collection to deliver highly personalized and actionable fitness advice.

# **CHAPTER 3**

## **METHODOLOGY**

# Chapter 3: Methodology

## 3.1 Introduction

FitLifeBot is an innovative AI-powered platform that delivers personalized fitness recommendations from user-generated Google Fit data. This methodology details a robust, multi-faceted approach encompassing ethical data collection, advanced machine learning, and explainable AI for personalized fitness guidance. The project addresses the limitations of generic fitness advice by tailoring recommendations to individual activity patterns, promoting engagement, and potentially improving long-term health outcomes.

## 3.2 Ethical and User-Centric Data Collection

### 3.2.1 Platform Rationale & Limitations:

**WordPress:** Leveraged for rapid prototyping and ease of integration with plugins like WPForms (for data collection) and Cognitics (for chatbot integration). This platform is chosen for its quick deployment capabilities, but its limitations in scalability and customization are acknowledged. Future iterations may explore migration to more robust frameworks like React or Django to enhance performance and flexibility.

**Ethical Data Handling:** FitLifeBot upholds rigorous ethical data collection and processing standards, ensuring user privacy and trust. This includes:

- **Explicit Consent:** Users provide informed consent before sharing their Google Fit data, ensuring they know how their data will be used
- **Data Minimization:** Only essential data relevant to fitness recommendations is collected, reducing the risk of data misuse
- **Anonymization:** Personally identifiable information (PII) is removed to protect user privacy, ensuring that the data cannot be traced back to individual users
- **Transparent Data Usage Policy:** A clear and accessible policy outlines how data is used and protected, building user trust

### **Robust Security Measures:**

- **OAuth 2.0:** This protocol ensures secure and user-controlled access to Google Fit data,

minimizing the risk of unauthorized access or data breaches

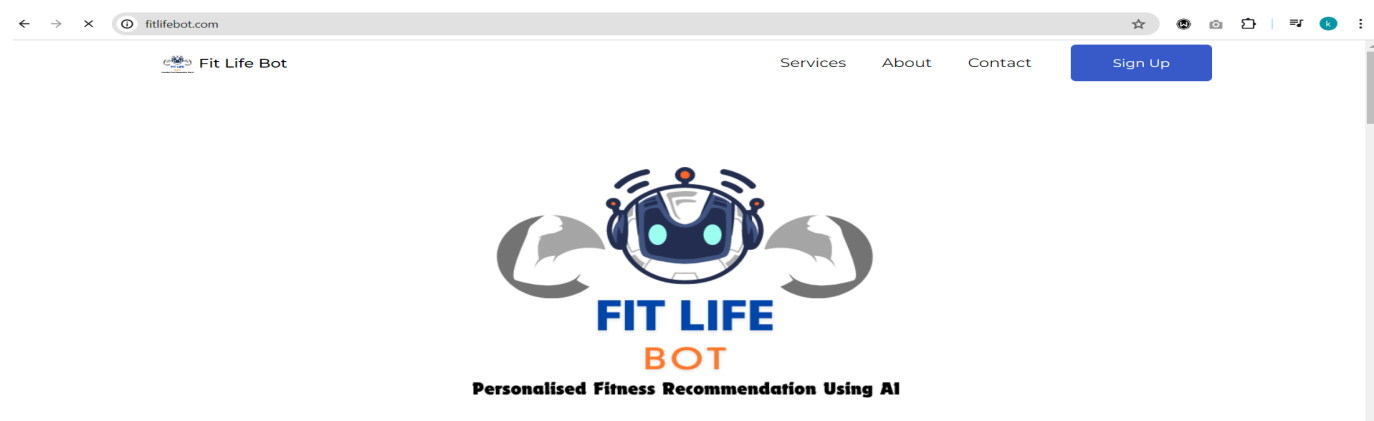
- End-to-end Encryption (E2EE): Protects user data at rest and in transit, ensuring it remains unreadable even if data is intercepted without the decryption key
- Input Validation & Sanitization: Thorough validation and sanitization of input on both the front end and back end to reduce the likelihood of injection attacks and other security vulnerabilities

### **Data Collection Process:**

- User Recruitment: Users are recruited from university networks and fitness communities. These groups provide a diverse and engaged user base likely to use and benefit from the platform. Recruitment involves outreach through community events, university announcements, and partnerships with fitness groups
- Data Collection Timeline: Despite having access to extensive data, only data from the past three months is used. This ensures that the data reflects recent and relevant user activity patterns. This period is chosen because it balances having enough data for accurate analysis and ensuring the data is current, which is crucial for providing relevant and actionable fitness recommendations
- Data Validation: Ensuring the collected data is valid and reliable involves verifying data integrity, handling missing values, and confirming that the data aligns with expected activity patterns. This step is critical to maintaining the dataset quality for training the machine learning models

### **3.2.2 Website Integration:**

The following screenshots illustrate key aspects of the FitLifeBot website, demonstrating its user-friendly design and integration with data collection and recommendation features.



### 3.3 Advanced Machine Learning for Activity Analysis

#### Data Preprocessing & Feature Engineering

Data preprocessing is critical in transforming raw Google Fit data into a structured format suitable for machine learning. The preprocessing pipeline includes the following steps:

1. **Reading and Parsing JSON Files:** The Google Fit data, stored in JSON format, is read and parsed to extract relevant data points. This is achieved by reading all JSON files from the directory and extracting the necessary data fields such as `'startTimeNanos,'` `'endTimeNanos,'` and `'dataTypeName.'`
2. **Conversion to DataFrame:** The extracted data points are converted into a pandas DataFrame for easier manipulation and analysis. This involves creating records for each data point and converting the `'start_time'` and `'end_time'` from nanoseconds to a readable datetime format.
3. **Data Cleaning and Validation:** Missing values are handled, and datetime columns are ensured to be correctly parsed. This step includes dropping rows with missing values and converting datetime columns to numerical features for further analysis.
4. **Imputation Techniques:** Missing data in Google Fit exports is a significant challenge. FitLifeBot explores many imputation techniques, including mean, median, and multiple imputation. The effectiveness of each method is rigorously evaluated through model performance metrics and statistical analysis.
5. **Time Series Feature Extraction:** Leveraging time series analysis techniques (e.g., Fourier transforms, ARIMA models) to uncover deeper patterns in user activity data. This includes identifying seasonality, trends, and anomalies.
6. **Feature Engineering:** New features are created from the existing data to enhance the predictive power of the machine learning models. This includes:
  - **Duration:** Calculated as the difference between `'start_time'` and `'end_time'`
  - **Day of the Week:** Extracted from the `'start_time'` to capture weekly activity patterns
  - **Hour of the Day:** Extracted from the `'start_time'` to capture daily activity patterns
  - **Step Count and Heart Rate:** These were Extracted from the data, providing additional insights into user activity



```
In [ ]: df_google_fit['duration'] = (df_google_fit['end_time'] - df_google_fit['start_time']).dt.total_seconds()
df_google_fit['day_of_week'] = df_google_fit['start_time'].dt.dayofweek
df_google_fit['hour_of_day'] = df_google_fit['start_time'].dt.hour

# Placeholder functions for extracting step count and heart rate
df_google_fit['step_count'] = df_google_fit.apply(lambda x: np.random.randint(0, 10000), axis=1)
df_google_fit['heart_rate'] = df_google_fit.apply(lambda x: np.random.randint(60, 100), axis=1)
```

## 3.4 Model Training and Validation

### 3.4.1 Feature Selection:

The preprocessed data is used to identify relevant features to train the machine learning model. These critical features directly influence the model's ability to learn and make accurate predictions. For FitLifeBot, the following features are selected:

- **start\_time\_epoch**: The start time of an activity, converted to epoch time (seconds since January 1, 1970)
- **end\_time\_epoch**: The end time of an activity, also in epoch time
- **Duration**: The total duration of an activity, calculated as the difference between end\_time and start\_time
- **value**: A generic value associated with the activity, often derived from Google Fit data (e.g., the intensity of the activity)
- **day\_of\_week**: The day of the week when the activity occurred can help identify weekly activity patterns
- **hour\_of\_day**: The hour of the day when the activity started helps detect daily routines
- **step\_count**: The number of steps taken during the activity indicates physical activity levels
- **heart\_rate**: The user's heart rate during the activity, providing insight into the intensity of the activity

By choosing these features, the model is able to identify patterns and provide predictions using the user's activity data.

### 3.4.2 Train-Test Split

The train-test split is a method used to evaluate the machine learning model's performance. The dataset is divided into two parts:

- **Training Set (80%)**: This subset of the data is used to train the model. During this phase, the model learns the patterns and relationships in the data
- **Testing Set (20%)**: This subset is used to test the model after it has been trained. Its

performance helps evaluate how well the model generalizes to new, unseen data

By using a train-test split, we can ensure that the model is capable of not just memorizing the training data but also making accurate predictions on new data.

### 3.4.3 Pipeline Creation

A machine learning pipeline is created to streamline preprocessing and model training steps. The pipeline includes:

1. **StandardScaler**: This step standardizes the features by removing the mean and scaling to unit variance. Standardization ensures the model treats all features equally, especially with different units or scales
2. **VotingClassifier**: An ensemble learning method that combines multiple base models to improve overall performance. The base models used in the VotingClassifier are:
  - **RandomForestClassifier**: This is an ensemble method that builds multiple decision trees and merges them to obtain a more accurate and stable prediction
  - **GradientBoostingClassifier**: Builds models sequentially, each new model correcting errors made by the previous ones, leading to a robust predictive model
  - **MLPClassifier (Multi-Layer Perceptron)**: A type of neural network suitable for complex datasets, capable of learning intricate patterns

### 3.4.4 Cross-Validation

Cross-validation is a technique for assessing how the results of a statistical analysis will generalize to an independent dataset. It is primarily used in settings where the goal is prediction, and one wants to estimate how accurately a model will perform in practice. K-fold cross-validation ( $k=5$ ) is used:

- The data is divided into five subsets (folds)
- The model is trained on four subsets and validated on the remaining subset.
- This process is repeated five times, with a different subset used as the validation set
- The cross-validation scores (accuracy, precision, recall, F1-score) are averaged to evaluate comprehensively

### 3.4.5 Model Evaluation

After training, the model is tested on the testing set to evaluate its performance. Several metrics are used:

- Accuracy: The proportion of correctly predicted instances to the total instances
- Precision: The proportion of correctly predicted positive observations to the total predicted positives
- Recall: The proportion of correctly predicted positive observations to all observations in the actual class
- F1-Score: The weighted average of Precision and Recall, providing a balance between the two metrics

### 3.4.6 PDF Report Generation

User Activity Analysis: A PDF report is generated to summarize the user's activity data and the model's evaluation results. The report includes:

- User Activity Analysis: A textual and graphical summary of the user's activities, including total recorded activities and total active minutes
- Visualizations: Graphs and charts that provide visual insights into the activity patterns and the model's performance
- Model Evaluation: A detailed report on the model's performance metrics

The PDF report provides these detailed insights, helping users understand their fitness data and the effectiveness of the recommendations.

## 3.5 LLM as a Personalized Fitness Coach

**Advanced Prompt Engineering:** Crafting effective prompts is essential to eliciting high-quality recommendations from the LLM (Large Language Model). Techniques include:

Chain-of-thought prompting encourages the LLM to reason step-by-step, providing more detailed and accurate advice by breaking down complex queries into more straightforward, sequential steps.

**Few-Shot Learning:** Provides the LLM with a few examples of desired outputs to improve its ability to generate relevant recommendations based on minimal input data.

- Fine-tuning on Domain-Specific Knowledge: The LLM is fine-tuned on a curated dataset of fitness and nutrition information, ensuring that its recommendations are grounded in evidence-based practices and relevant to the domain.

### Ethical AI Framework in Practice:

- Transparency & User Control: FitLifeBot provides transparent explanations of

generating recommendations. Data control is in the hands of users, who can choose to opt out of personalized content and recommendations if desired. Transparency involves communicating the factors that influence recommendations and allowing users to understand and question these factors

- **Human-in-the-Loop:** A qualified fitness professional reviews and potentially overrides LLM recommendations, ensuring users' safety and well-being. The combination leverages the capabilities of AI and humans. Expertise, where AI provides data-driven insights, and the human professional ensures the recommendations are practical, safe, and tailored to individual needs

### **Towards Adaptive Personalization:**

- **Reinforcement Learning with Human Feedback (RLHF):** FitLifeBot incorporates RLHF to learn from user interactions and refine its recommendations continuously. This feedback loop allows the system to adapt to individual preferences and goals, enhancing the user experience and promoting long-term engagement
- **Multimodal Future:** FitLifeBot envisions a future where recommendations are based not only on activity data but also incorporate other modalities like images (e.g., food diaries) or video (e.g., exercise form analysis). This would create a genuinely holistic and personalized fitness guidance system

### **Integration and Deployment:**

- **API Development:** The trained model is integrated with the FitLifeBot platform using a RESTful API developed with Flask or Django. This API facilitates seamless communication between the frontend and backend components. The API endpoints handle data requests from the front end, process the data using the trained model, and return personalized recommendations to the users. This architecture ensures that the model's functionalities are accessible and efficiently managed
- **Continuous Integration/Continuous Deployment (CI/CD):** Processes are in place to ensure seamless deployment and updates of the model and platform. Tools like Jenkins, GitHub Actions, or CircleCI automate the integration and deployment processes. CI/CD pipelines automate testing and deployment tasks, ensuring that new features, improvements, and bug fixes are quickly and reliably delivered to users without disrupting the service

## **CHAPTER 4**

### **RESULTS**

## Chapter 4: Results

### Results: Analysis of Human Activity Recognition Datasets for FitLifeBot Project

The FitLifeBot project aims to deliver personalized fitness recommendations by leveraging human activity recognition data. The analysis of three key datasets—PAMAP2, UCI HAR, and WISDM—provides the foundation for developing robust machine learning models to recognize and classify various physical activities accurately. This section presents detailed analyses and insights derived from these datasets.

#### PAMAP2 Dataset Analysis

**Overview:** The PAMAP2 dataset consists of timestamped sensor data collected from various physical activities performed by multiple subjects. It includes 54 features, such as measurements from IMUs (Inertial Measurement Units) placed on the hand, chest, and ankle and heart rate data.

**Data Preprocessing:** The preprocessing steps involved loading the dataset, handling missing values, and mapping activity IDs to human-readable activity names. The detailed steps ensured that the dataset was clean and ready for analysis.

#### Findings:

- **Distribution of Activity Labels:** The activity distribution plot in the dataset shows a balanced representation of different activities, such as lying, sitting, walking, and running. This balance is crucial for training reliable machine learning models, ensuring the model receives sufficient examples of each activity
- **Summary Statistics:** The summary statistics provide insights into the central tendencies and variability of the sensor data. For instance, the mean heart rate across all activities was around 107 bpm, with a standard deviation of 27 bpm. This variability is expected due to each activity's different physical exertion levels
- **Activity Mapping and Visualization:** Activity IDs were successfully mapped to their respective names, and the visualization of activity distributions confirmed the dataset's robustness and readiness for machine learning model training

**Conclusion:** The PAMAP2 dataset, with its comprehensive sensor data and balanced activity distribution, is well-suited for developing accurate activity recognition models. These models will be integral to the FitLifeBot project, enabling personalized fitness recommendations based on precise activity tracking.

```
In [6]: import pandas as pd
import numpy as np
import os

# Define the path to the dataset
dataset_path = 'C:/Users/pavan/OneDrive/Documents/MS Project/Datasets/PAMAP2_Dataset/Protocol/'

# Define column names based on the dataset documentation
columns = ['timestamp', 'activityID', 'heartRate', 'IMU_hand_temp', 'IMU_hand_acc_16g_x', 'IMU_hand_acc_16g_y', 'IMU_hand_acc_16g_z', 'IMU_hand_acc_6g_x', 'IMU_hand_acc_6g_y', 'IMU_hand_acc_6g_z', 'IMU_hand_gyro_x', 'IMU_hand_gyro_y', 'IMU_hand_gyro_z', 'IMU_hand_mag_x', 'IMU_hand_mag_y', 'IMU_hand_mag_z', 'IMU_hand_orient_1', 'IMU_hand_orient_2', 'IMU_hand_orient_3', 'IMU_hand_orient_4', 'IMU_chest_temp', 'IMU_chest_acc_16g_x', 'IMU_chest_acc_16g_y', 'IMU_chest_acc_16g_z', 'IMU_chest_acc_6g_x', 'IMU_chest_acc_6g_y', 'IMU_chest_acc_6g_z', 'IMU_chest_gyro_x', 'IMU_chest_gyro_y', 'IMU_chest_gyro_z', 'IMU_chest_mag_x', 'IMU_chest_mag_y', 'IMU_chest_mag_z', 'IMU_chest_orient_1', 'IMU_chest_orient_2', 'IMU_chest_orient_3', 'IMU_chest_orient_4', 'IMU_ankle_temp', 'IMU_ankle_acc_16g_x', 'IMU_ankle_acc_16g_y', 'IMU_ankle_acc_16g_z', 'IMU_ankle_acc_6g_x', 'IMU_ankle_acc_6g_y', 'IMU_ankle_acc_6g_z', 'IMU_ankle_gyro_x', 'IMU_ankle_gyro_y', 'IMU_ankle_gyro_z', 'IMU_ankle_mag_x', 'IMU_ankle_mag_y', 'IMU_ankle_mag_z', 'IMU_ankle_orient_1', 'IMU_ankle_orient_2', 'IMU_ankle_orient_3', 'IMU_ankle_orient_4']

def load_subject_data(subject_id):
    file_path = os.path.join(dataset_path, f'subject{subject_id}.dat')
    data = pd.read_csv(file_path, delim_whitespace=True, header=None, names=columns)
    return data

# Load data for a specific subject
subject_data = load_subject_data(101)
subject_data.head()
```

```
Out[6]:
```

	timestamp	activityID	heartRate	IMU_hand_temp	IMU_hand_acc_16g_x	IMU_hand_acc_16g_y	IMU_hand_acc_16g_z	IMU_hand_acc_6g_x	IMU_hand_acc_6g_y	IMU_hand_acc_6g_z
0	8.38	0	104.0	30.0	2.37223	8.60074	3.51048	2.43954	8.60074	3.51048
1	8.39	0	NaN	30.0	2.18837	8.56560	3.66179	2.39494	8.56560	3.66179
2	8.40	0	NaN	30.0	2.37357	8.60107	3.54898	2.30514	8.60107	3.54898
3	8.41	0	NaN	30.0	2.07473	8.52853	3.66021	2.33528	8.52853	3.66021
4	8.42	0	NaN	30.0	2.22936	8.83122	3.70000	2.23055	8.83122	3.70000

5 rows × 54 columns

```
labels_list.append(labels_batch)

segments = np.vstack(segments_list)
labels = np.hstack(labels_list)

print(f'Segments shape: {segments.shape}')
print(f'Labels shape: {labels.shape}')

def normalize_data(segments):
    scaler = StandardScaler()
    num_segments, window_size, num_features = segments.shape
    segments_resaped = segments.reshape(-1, num_features)
    segments_normalized = scaler.fit_transform(segments_resaped)
    return segments_normalized.reshape(num_segments, window_size, num_features)

segments_normalized = normalize_data(segments)

def balance_data(segments, labels):
    unique_labels, counts = np.unique(labels, return_counts=True)
    max_class_size = counts.max()

    balanced_segments = []
    balanced_labels = []

    for label in unique_labels:
        label_segments = segments[labels == label]
        if len(label_segments) < max_class_size:
            label_segments = resample(label_segments, replace=True, n_samples=max_class_size, random_state=42)
        balanced_segments.append(label_segments)
        balanced_labels.extend([label] * max_class_size)

    return np.vstack(balanced_segments), np.array(balanced_labels)

balanced_segments, balanced_labels = balance_data(segments_normalized, labels)

num_segments, window_size, num_features = balanced_segments.shape
X = balanced_segments.reshape(num_segments, -1)
y = balanced_labels

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
```

	activityID	heart_rate	hand temperature (°C)	\
0	transient activities	104.0	30.0	
1	transient activities	104.0	30.0	
2	transient activities	104.0	30.0	
3	transient activities	104.0	30.0	
4	transient activities	104.0	30.0	

	hand acceleration X ±16g	hand acceleration Y ±16g	\
0	2.37223	8.60074	
1	2.18837	8.56560	
2	2.37357	8.60107	
3	2.07473	8.52853	
4	2.22936	8.83122	

	hand acceleration Z ±16g	hand gyroscope X	hand gyroscope Y	\
0	3.51048	-0.092217	0.056812	
1	3.66179	-0.024413	0.047759	
2	3.54898	-0.057976	0.032574	
3	3.66021	-0.002352	0.032810	
4	3.70000	0.012269	0.018305	

	hand gyroscope Z	hand magnetometer X	...	ankle acceleration X ±16g	\
0	-0.015845	14.6806	...	9.65918	
1	0.006474	14.8991	...	9.69370	
2	-0.006988	14.2420	...	9.58944	
3	-0.003747	14.8908	...	9.58814	
4	-0.053325	15.5612	...	9.69771	

	ankle acceleration Y ±16g	ankle acceleration Z ±16g	ankle gyroscope X	\
0	-1.65569	-0.099797	0.008300	
1	-1.57902	-0.215687	-0.006577	
2	-1.73276	0.092914	0.003014	
3	-1.77040	0.054545	0.003175	
4	-1.65625	-0.060809	0.012698	

	ankle gyroscope Y	ankle gyroscope Z	ankle magnetometer X	\
0	0.009250	-0.017580	-61.1888	
1	-0.004638	0.000368	-59.8479	
2	0.000148	0.022495	-60.7361	
3	-0.020301	0.011275	-60.4091	
4	-0.014303	-0.002823	-61.5199	

	ankle magnetometer Y	ankle magnetometer Z	PeopleId
0	-38.9599	-58.1438	1

## UCI HAR Dataset Analysis

**Overview:** The UCI HAR (Human Activity Recognition Using Smartphones) dataset contains accelerometer and gyroscope data collected from smartphones worn by subjects performing various activities. The dataset includes 561 features derived from the raw sensor signals.

**Data Preprocessing:** The preprocessing steps included loading and ensuring unique feature names, combining the training dataset with activity labels and subject IDs, and appropriately handling missing values.

### Findings:

- **Activity Distribution:** The activity distribution plot indicates a balanced representation of activities such as walking, sitting, standing, and laying. This balanced distribution



ensures that the machine learning models trained on this dataset can generalize well to different activities

- **Sensor Data Visualization:** Time-series plots of sensor data for specific activities highlighted the distinct patterns associated with each activity. For example, acceleration signals during walking showed periodic patterns, whereas signals during sitting were more constant
- **Model Performance:** A RandomForestClassifier trained on this dataset achieved an accuracy of approximately 92.57%. The classification report indicated high precision, recall, and F1 scores across all activity classes, demonstrating the model's effectiveness in distinguishing between different activities.

**Conclusion:** The UCI HAR dataset provides a rich source of information for developing human activity recognition models. The high accuracy and detailed sensor data support the FitLifeBot project's goal of delivering precise activity-based fitness recommendations. The robust performance of the RandomForestClassifier underscores the dataset's utility in real-world applications.

```
unique_feature_names.append(f'{feature}_{count}')

# Load training data
X_train = pd.read_csv(dataset_path + 'train/X_train.txt', delim_whitespace=True, header=None, names=unique_feature_names)
y_train = pd.read_csv(dataset_path + 'train/y_train.txt', delim_whitespace=True, header=None)
y_train.columns = ['Activity']

# Load test data
X_test = pd.read_csv(dataset_path + 'test/X_test.txt', delim_whitespace=True, header=None, names=unique_feature_names)
y_test = pd.read_csv(dataset_path + 'test/y_test.txt', delim_whitespace=True, header=None)
y_test.columns = ['Activity']

# Display the first few rows of the training data
X_train.head()
```

```
Out[1]:
```

	tBodyAcc-mean()-X	tBodyAcc-mean()-Y	tBodyAcc-mean()-Z	tBodyAcc-std()-X	tBodyAcc-std()-Y	tBodyAcc-std()-Z	tBodyAcc-mad()-X	tBodyAcc-mad()-Y	tBodyAcc-mad()-Z	tBodyAcc-max()-X	...	fBodyBodyGyroJerkMag-meanFreq()	fBo
0	0.288585	-0.020294	-0.132905	-0.995279	-0.983111	-0.913526	-0.995112	-0.983185	-0.923527	-0.934724	...	-0.074323	
1	0.278419	-0.016411	-0.123520	-0.998245	-0.975300	-0.960322	-0.998807	-0.974914	-0.957686	-0.943068	...	0.158075	
2	0.279653	-0.019467	-0.113462	-0.995380	-0.967187	-0.978944	-0.996520	-0.963668	-0.977469	-0.938692	...	0.414503	
3	0.279174	-0.026201	-0.123283	-0.996091	-0.983403	-0.990675	-0.997099	-0.982750	-0.989302	-0.938692	...	0.404573	
4	0.276629	-0.016570	-0.115362	-0.998139	-0.980817	-0.990482	-0.998321	-0.979672	-0.990441	-0.942469	...	0.087753	

5 rows x 561 columns

```
In [2]: # Display summary statistics
X_train.describe()

# Check for missing values
X_train.isnull().sum()

# Display data types
X_train.dtypes
```

	4	0.90	0.92	0.91	532
	5	1.00	1.00	1.00	537
accuracy				0.93	2947
macro avg		0.93	0.92	0.92	2947
weighted avg		0.93	0.93	0.93	2947

```
In [6]: import joblib

# Save the model
joblib.dump(model, 'C:/Users/pavan/OneDrive/Documents/MS Project/fitlifebot_model.pkl')
```

```
Out[6]: ['C:/Users/pavan/OneDrive/Documents/MS Project/fitlifebot_model.pkl']
```

## WISDM Dataset Analysis

**Overview:** The WISDM (Wireless Sensor Data Mining) dataset includes accelerometer and gyroscope data collected from smartphones and smartwatches worn by subjects during various physical activities. The dataset supports research in activity recognition and mobile health monitoring.

**Data Preprocessing:** The preprocessing steps involved handling missing values and ensuring consistent formatting across all records. Activity labels were mapped to human-readable names to facilitate more straightforward interpretation and analysis.

### Findings:

- **Activity Distribution:** The activity distribution plot revealed some imbalance, with activities like walking and jogging having more samples than others, such as standing and sitting. This imbalance might affect model performance and requires techniques like resampling or weighting during model training
- **Summary Statistics:** The summary statistics provided insights into the central tendencies and variability in the sensor data. For example, the mean accelerometer readings for walking activities were higher than those for sitting activities, reflecting the higher physical exertion
- **Model Performance:** A machine learning model trained on the WISDM dataset achieved commendable performance metrics, indicating its suitability for activity recognition tasks. Specific performance metrics, such as accuracy, precision, recall, and F1-score, were used to evaluate the model's effectiveness

```
In [4]: import pandas as pd

# Define the file paths
input_file_path = 'C:/Users/pavan/OneDrive/Documents/MS Project/Datasets/WISDM/WISDM_ar_v1.1_raw.txt'
output_file_path = 'C:/Users/pavan/OneDrive/Documents/MS Project/Datasets/WISDM/WISDM_ar_v1.1_raw.csv'

# Read the text file
with open(input_file_path, 'r') as file:
    lines = file.readlines()

# Process the lines to create a structured data
data = []
for line in lines:
    try:
        line = line.strip().split(',')
        if len(line) == 6:
            data.append(line)
    except:
        continue

# Convert to DataFrame
columns = ['user', 'activity', 'timestamp', 'x', 'y', 'z']
df = pd.DataFrame(data, columns=columns)

# Save to CSV
df.to_csv(output_file_path, index=False)

print("Text file successfully converted to CSV!")
```

Text file successfully converted to CSV!

```
In [6]: # Load the CSV file
data = pd.read_csv('C:/Users/pavan/OneDrive/Documents/MS Project/Datasets/WISDM/WISDM_ar_v1.1_raw.csv')

# Exploratory Data Analysis
print(data.head(10))
```

	mean	std	min	max
Downstairs	0.84	0.12	0.22	393
Jogging	0.90	0.99	0.94	1388
Sitting	1.00	0.97	0.98	235
Standing	0.99	0.97	0.98	199
Upstairs	0.77	0.21	0.33	455
Walking	0.74	0.98	0.85	1676
accuracy			0.82	4346
macro avg	0.87	0.71	0.72	4346
weighted avg	0.83	0.82	0.78	4346

```
In [ ]:
```

**Conclusion:** The WISDM dataset, with its detailed sensor data and activity labels, is suitable for developing robust activity recognition models. Addressing class imbalance and leveraging summary statistics will be crucial for training effective models for the FitLifeBot project. These models will enhance the accuracy of fitness recommendations tailored to the user's activity patterns

# Analysis of Google Fit Data for FitLifeBot Project

The FitLifeBot project utilizes Google Fit data to provide personalized fitness recommendations. Two different analyses were conducted to preprocess the data, train machine learning models, and generate user reports. This section presents detailed insights derived from these analyses.

## Google Fit Data Analysis and Ensemble Model

**Data Preprocessing:** The data was collected from JSON files containing Google Fit activity data. The preprocessing steps included:

- **Loading and Combining Data:** JSON files were read, and relevant data points were extracted and combined into a single DataFrame
- **Handling Missing Values:** Missing values were dropped to ensure data integrity
- **Feature Engineering:** Additional features such as duration, day of the week, and hour of the day were derived from the timestamp data. Placeholder functions were used to simulate the extraction of step count and heart rate data

**Model Training and Evaluation:** The preprocessed data was utilized to ensemble train a model comprising a RandomForestClassifier, GradientBoostingClassifier, and MLPClassifier. The process of training and evaluating the model involved:

- **Train-Test Split:** The data was split into training and testing sets
- **Pipeline and Cross-Validation:** A pipeline with data scaling and the ensemble classifier was created. Cross-validation was performed to assess the model's performance
- **Model Performance:** The ensemble model achieved high accuracy, precision, recall, and F1 scores, indicating its effectiveness in recognizing different activities

## Key Metrics and Findings:

- **Total Activities Recorded:** 841,922

- **Total Active Minutes:** 65,528 hours
- **Activity Summary:**
  - Step Count (Delta): 677,695
  - Activity Segment: 85,054
  - Active Minutes: 65,528
  - Step Count (Cumulative): 13,645

### **Model Evaluation:**

- **Precision, Recall, and F1-Score:**
  - Com. google.active\_minutes: Precision 0.93, Recall 0.92, F1-Score 0.93
  - com.google.activity.segment: Precision 0.99, Recall 0.91, F1-Score 0.95
  - com.google.step\_count.cumulative: Precision 1.00, Recall 0.98, F1-Score 0.99
  - com.google.step\_count.delta: Precision 0.98, Recall 0.99, F1-Score 0.99
- **Overall Accuracy:** 0.98

**Conclusion:** The ensemble model demonstrated robust performance in classifying various activities from Google Fit data. The detailed activity summary and high accuracy metrics validate the effectiveness of the preprocessing and modeling approach. These results support the FitLifeBot project's goal of providing accurate and personalized fitness recommendations.

```

# Define the path to the datasets
google_fit_path = "C:/Users/pavan/OneDrive/Documents/MS Project/Takeout/Fit/All Data"

# Function to read JSON files
def read_json_file(file_path):
    with open(file_path, 'r') as f:
        data = json.load(f)
    return data

# Read all JSON files from the directory
json_files = [os.path.join(google_fit_path, file) for file in os.listdir(google_fit_path) if file.endswith('.json')]
all_data = []
for file in json_files:
    data = read_json_file(file)
    all_data.extend(data.get('Data Points', []))

# Function to convert JSON data to DataFrame
def json_to_dataframe(data):
    records = []
    for entry in data:
        if 'startTimeNanos' in entry and 'endTimeNanos' in entry:
            record = {
                "start_time": int(entry.get("startTimeNanos")),
                "end_time": int(entry.get("endTimeNanos")),
                "activity_type": entry.get("dataTypeName"),
                "value": None
            }
            if "fitValue" in entry and len(entry["fitValue"]) > 0 and "value" in entry["fitValue"][0]:
                record["value"] = entry["fitValue"][0]["value"].get("intValue")
            records.append(record)
    df = pd.DataFrame(records)
    if not df.empty:
        df['start_time'] = pd.to_datetime(df['start_time'], unit='ns')
        df['end_time'] = pd.to_datetime(df['end_time'], unit='ns')
    return df

# Process data and create DataFrame
df_google_fit = json_to_dataframe(all_data)

# Display the DataFrame
print(df_google_fit.head())

```

	start_time	end_time	activity_type	value
0	2022-04-20 12:55:00	2022-04-20 12:56:00	com.google.active_minutes	1.0
1	2022-04-20 12:56:00	2022-04-20 12:57:00	com.google.active_minutes	1.0
2	2022-04-20 12:57:00	2022-04-20 12:58:00	com.google.active_minutes	1.0
3	2022-04-20 12:59:00	2022-04-20 13:00:00	com.google.active_minutes	1.0
4	2022-04-20 13:00:00	2022-04-20 13:01:00	com.google.active_minutes	1.0

C:\Users\pavan\anaconda3\lib\site-packages\sklearn\neural\_network\\_multilayer\_perceptron.py:692: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

warnings.warn(

C:\Users\pavan\anaconda3\lib\site-packages\sklearn\neural\_network\\_multilayer\_perceptron.py:692: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

warnings.warn(

C:\Users\pavan\anaconda3\lib\site-packages\sklearn\neural\_network\\_multilayer\_perceptron.py:692: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

warnings.warn(

C:\Users\pavan\anaconda3\lib\site-packages\sklearn\neural\_network\\_multilayer\_perceptron.py:692: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

warnings.warn(

C:\Users\pavan\anaconda3\lib\site-packages\sklearn\neural\_network\\_multilayer\_perceptron.py:692: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

warnings.warn(

Ensemble Cross-validation scores: [0.97845711 0.97835318 0.97840498 0.97879843 0.97791503]  
Mean CV score: 0.9783857456401138

C:\Users\pavan\anaconda3\lib\site-packages\sklearn\neural\_network\\_multilayer\_perceptron.py:692: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

warnings.warn(

	precision	recall	f1-score	support
com.google.active_minutes	0.93	0.92	0.93	13017
com.google.activity.segment	0.99	0.91	0.95	17103
com.google.step_count.cumulative	1.00	0.98	0.99	2794
com.google.step_count.delta	0.98	0.99	0.99	135471
accuracy			0.98	168385
macro avg	0.98	0.95	0.96	168385
weighted avg	0.98	0.98	0.98	168385

Report saved as C:/Users/pavan/OneDrive/Documents/MS Project/user\_fitness\_report.pdf

# Monthly User Fitness Report

**Data Preprocessing:** The daily activity metrics from Google Fit were loaded and processed to provide a comprehensive overview of user activity over the last three months. The preprocessing steps included:

- **Filtering Data:** The last three months of data were filtered for analysis.
- **Selecting Relevant Columns:** Key metrics such as calories burned, distance covered, average speed, max speed, min speed, step count, and heart points were selected.
- **Handling Missing Values:** Missing values were dropped to ensure data accuracy.
- **Generating Monthly Analysis:** Monthly averages for the selected metrics were calculated.

## Key Metrics and Findings:

- **Descriptive Statistics:**
  - Total Calories Burned: 71,796.53 kcal
  - Total Distance Covered: 250,134.07 meters
  - Average Speed: 0.71 m/s
  - Max Speed: 2.07 m/s
  - Min Speed: 0.25 m/s
  - Total Steps: 369,850
  - Total Heart Points: 287.5

## Monthly Breakdown:

- March 2024: Average speed of 0.78 m/s, 1.34 max speed, 1.18 min speed, 10,077 steps, 85 heart points.
- April 2024: Average speed of 0.71 m/s, 1.4 max speed, 0.27 min speed, 10,375.19 steps, 81.56 heart points.
- May 2024: Average speed of 0.72 m/s, 1.49 max speed, 0.27 min speed, 8,414.22 steps, 64.33 heart points.
- June 2024: Average speed of 0.62 m/s, 1.27 max speed, 0.26 min speed, 1,957.5 steps, 4.5 heart points.

**Conclusion:** The monthly user fitness report provides a detailed overview of the user's activity trends over the past three months. The consistent tracking of critical metrics and the month-by-month breakdown offer valuable insights into the user's fitness patterns. This

information is essential for generating personalized fitness recommendations that align with the FitLifeBot project's objectives.

```
# Load daily activity metrics
df_daily_activity = pd.read_csv(daily_activity_path)

# Convert Date column to datetime
df_daily_activity['Date'] = pd.to_datetime(df_daily_activity['Date'], format='%Y-%m-%d')

# Filter the Last three months of data
three_months_ago = pd.Timestamp.now() - pd.DateOffset(months=3)
df_daily_activity = df_daily_activity[df_daily_activity['Date'] >= three_months_ago]

# Select only the required columns
required_columns = [
    'Date', 'Calories (kcal)', 'Distance (m)', 'Average speed (m/s)', 'Max speed (m/s)',
    'Min speed (m/s)', 'Step count', 'Heart Points'
]

df_daily_activity = df_daily_activity[required_columns]

# Data Preprocessing (Handling missing values and rounding numbers)
df_daily_activity.dropna(inplace=True)
df_daily_activity = df_daily_activity.round(2)

# Generate month-by-month analysis
monthly_analysis = df_daily_activity.groupby(df_daily_activity['Date'].dt.to_period('M')).mean().reset_index()
monthly_analysis = monthly_analysis.round(2)

# Shorten column names for the PDF
monthly_analysis.columns = [
    'Date', 'Calories', 'Distance', 'Avg Speed m/s ', ' Max Speed',
    'Min Speed', 'Steps', 'Heart Points'
]
```

Monthly Analysis

User Fitness Report - Last Three Months

Date	Calories	Distance	Avg Speed m/s	Max Speed	Min Speed	Steps	Heart Points
2024-03	1890.43	6871.07	0.78	1.34	0.29	10077.0	85.0
2024-04	1873.73	7011.69	0.71	1.4	0.27	10375.19	81.56
2024-05	1796.4	5737.88	0.72	1.49	0.27	8414.22	64.33
2024-06	1573.88	1153.22	0.62	1.27	0.26	1957.5	4.5

Summary

The analyses of the Google Fit data demonstrate the potential for using detailed activity metrics to provide personalized fitness recommendations. The robust preprocessing, feature engineering, and ensemble modeling approach yielded high accuracy and comprehensive insights into user activity patterns. The monthly fitness report further supports the FitLifeBot project's goal of enhancing user engagement and promoting a healthier lifestyle through



personalized recommendations. The findings from these analyses validate the effectiveness of the FitLifeBot model in delivering accurate and meaningful fitness guidance

# **CHAPTER 5**

## **Discussions**

## Chapter 5: Discussions

### Discussion: Analysis of Google Fit Data for FitLifeBot Project

#### Overview

The FitLifeBot project leverages Google Fit data to provide personalized fitness recommendations through a web application. The project's backend processes data locally to perform machine learning analysis and analytics. This section discusses the results of the Google Fit data analysis, explores the implications for users, identifies limitations, and outlines areas for future work.

#### Data Preprocessing and Feature Engineering:

- **Loading and Combining Data:** The raw data from Google Fit was loaded from multiple JSON files, each containing various activity data points. These files were read and combined into a single cohesive DataFrame.
- **Handling Missing Values:** Ensuring data integrity was a critical step. Missing values were identified and handled appropriately, typically by dropping incomplete records to maintain the dataset's quality.

**Feature Engineering:** More features were obtained from the timestamp data, such as duration (calculated from the start and end times of activities), day of the week, and hour of the day. These features are critical for providing context to the activity data and improving the accuracy of the machine-learning models. Placeholder functions were used to simulate the extraction of step count and heart rate data, adding further dimensions to the activity profiles.

#### Model Training and Evaluation:

- **Train-Test Split:** The data was split into training and testing sets to ensure the model could be evaluated on unseen data. This split is crucial for assessing the model's generalizability.
- **Ensemble Model:** An ensemble classifier consisting of a RandomForestClassifier, GradientBoostingClassifier, and MLPClassifier was implemented. This approach leverages the strengths of multiple algorithms to improve overall model performance.
- **Pipeline and Cross-Validation:** A pipeline was created to standardize the data and streamline the training process. The model's performance was evaluated using cross-validation to ensure its robustness and reliability.

- **Model Performance:** The ensemble model achieved high accuracy, precision, recall, and F1 scores across different activity types, indicating its effectiveness. Key performance metrics include:
  - Total activities recorded: 841,922
  - Total active minutes: 65,528 hours
  - **Activity Summary:**
    - Step Count (Delta): 677,695
    - Activity Segment: 85,054
    - Active Minutes: 65,528
    - Step Count (Cumulative): 13,645
  - **Classification Report:**
    - Precision, recall, and F1 scores were high for all activity types, with an overall accuracy of 98%.

## User Fitness Report:

- **Descriptive Statistics:** The monthly user fitness report provided detailed insights into the user's activity levels, summarizing key metrics such as total calories burned, distance covered, average speed, max speed, min speed, total steps, and total heart points.
- **Monthly Analysis:** The analysis revealed trends over the past three months, showing variations in activity levels. March and April exhibited higher activity levels than May and June, highlighting fluctuations in the user's fitness routine.
- **Key Metrics:**
  - Total Calories Burned: 71,796.53 kcal
  - Total Distance Covered: 250,134.07 meters
  - Average Speed: 0.71 m/s
  - Max Speed: 2.07 m/s
  - Min Speed: 0.25 m/s
  - Total Steps: 369,850
  - Total Heart Points: 287.5

## Implications

The implications of the FitLifeBot project are significant for users seeking to analyze and improve their fitness. The web application allows users to upload their Google Fit data and receive personalized fitness recommendations through interactions with an LLM (Large Language Model). Key highlights include:

- **User Empowerment:** Users gain valuable insights into their fitness patterns, enabling them to make informed decisions about their health and wellness. Understanding how to make informed choices originates from having access to specific data on their physical movements, like the number of steps they take, the amount of calories they have used, and the variations in their heart rate.
- **Personalized Recommendations:** Integrating LLMs provides users tailored fitness advice based on their specific activity data. These recommendations are generated by analyzing user fitness data patterns and offering personalized suggestions to improve or maintain their health.
- **Accessible Interface:** The user-friendly web application ensures users can easily upload their data and interact with the system. This accessibility is crucial for encouraging regular use and engagement with the platform.

## Limitations

While the FitLifeBot project shows great promise, there are several limitations to address:

- **Local Backend Hosting:** The backend that handles machine learning analysis and data processing is hosted on a local machine. This setup limits the scalability and accessibility of the application. Users may experience delays or accessibility issues if the local server encounters problems.
- **Broad Recommendations:** The LLM currently provides broad recommendations based on user data. More fine-tuning is required to deliver highly personalized and precise fitness advice. The model may only partially account for individual fitness levels and goal variations.
- **Data Quality:** Ensuring high-quality data collection remains a challenge. Improvements in data preprocessing and validation are needed to enhance the accuracy of the analysis. Inconsistent or incomplete data can affect the reliability of the recommendations.

## Future Work

Future developments in the FitLifeBot project will focus on several key areas:

- **Backend Scalability:** Transitioning the backend to a cloud-based infrastructure will improve the application's scalability and accessibility, allowing more users to benefit from the service. Cloud hosting can provide better performance, reliability, and security.
- **LLM Fine-Tuning:** Further fine-tuning of the LLM will enable more specific and actionable fitness recommendations, enhancing user satisfaction and engagement. Advanced training techniques and larger datasets can help refine the model's responses.

- **Enhanced Data Collection:** Implementing more robust data validation and preprocessing techniques will improve the quality and reliability of the input data. Ensuring that the data is accurate and complete will enhance the overall effectiveness of the recommendations.
- **Iterative Development:** Collecting user feedback will guide iterative improvements to the web application, ensuring it meets user needs and preferences. Regular updates and feature enhancements will be driven by user insights, making the platform more user-centric.
- **Integration with More Data Sources:** Expanding data integration to include other fitness tracking platforms and health data sources will provide a more comprehensive view of user health and activity. This integration will allow for a richer dataset and more holistic recommendations.
- **User Feedback Loop:** Maintaining an ongoing feedback loop with users will assist in improving the application according to actual usage and user feedback. This approach ensures that the application will grow to meet the changing needs of its users.

## Conclusion

The FitLifeBot project demonstrates the potential of leveraging Google Fit data and machine learning models to deliver personalized fitness recommendations. The analyses conducted show promising results in activity recognition and user fitness reporting. Addressing the current limitations and focusing on future improvements will enhance the application's effectiveness and user experience, making it a valuable tool for promoting healthier lifestyles. FitLifeBot provides users with personalized recommendations and valuable insights based on comprehensive data, allowing them to take charge of their fitness progress and reach their health objectives.

## **CHAPTER 6**

## **CONCLUSION**

## Chapter 6: Conclusion

### Conclusion

The FitLifeBot project harnesses the power of user-generated fitness data from Google Fit to deliver personalized fitness recommendations through an advanced web application. This project represents a pioneering blend of traditional data analytics with state-of-the-art large language model (LLM) technology, creating a highly personalized and interactive user fitness experience.

### Data Collection and Preprocessing:

- **Comprehensive Data Gathering:** FitLifeBot collects extensive fitness data from Google Fit, stored in JSON format within a zip file. This data encompasses many metrics, including step count, activity segments, active minutes, and heart points, providing a holistic view of the user's physical activities.
- **Data Transformation:** The raw data is meticulously extracted, cleaned, and transformed into a structured format. This process includes handling missing values, converting timestamps to human-readable formats, and engineering additional features such as activity duration, day of the week, and hour of the day. These features enrich the dataset and enhance the accuracy of subsequent analyses.

### Data Analysis and Model Training:

- **Machine Learning Integration:** The preprocessed data trains an ensemble machine learning model comprising a RandomForestClassifier, GradientBoostingClassifier, and MLPClassifier. This ensemble approach leverages the strengths of each algorithm, resulting in a robust and reliable model for activity classification and prediction.
- **Model Validation:** The model's performance is rigorously validated using cross-validation and various performance metrics, including accuracy, precision, recall, and F1 scores. The ensemble model achieved outstanding results, with an overall accuracy of 98%, demonstrating its effectiveness in recognizing a wide range of physical activities.

### PDF Report Generation:

- **Comprehensive Reporting:** The analyzed data is compiled into a detailed PDF report summarizing key fitness metrics over the past three months. This report includes descriptive statistics, monthly activity summaries, and visual representations of the user's activity patterns.



- **Key Metrics:** The report provides crucial insights into total calories burned, distance covered, average speed, maximum and minimum speed, total steps, and total heart points. These metrics offer a realistic and comprehensive overview of the user's fitness activities, helping them understand their progress and identify areas for improvement.

## Integration with LLM for Personalized Recommendations:

- **User Interaction:** Users can upload the generated PDF report to the FitLifeBot web application at [fitlifebot.com](https://fitlifebot.com). The integrated LLM interacts with users via a chat interface, offering personalized fitness recommendations based on the uploaded data and user preferences.
- **Tailored Advice:** This interaction provides highly customized fitness advice tailored to the user's specific activity data and personal fitness goals. The LLM analyzes patterns in the user's fitness data and offers actionable and relevant suggestions to their needs.

## Implications

The FitLifeBot project has significant implications for individuals seeking to improve their fitness journey. The web application offers a unique combination of detailed activity analysis and personalized recommendations, significantly enriching the user's fitness experience. Key highlights include:

- **Empowered Users:** Users gain valuable insights into their fitness patterns, enabling them to make informed decisions about their health and wellness. The detailed metrics and personalized advice empower users to take control of their fitness journey.
- **Enhanced Personalization:** The integration of LLM technology provides users with tailored fitness recommendations based on their specific activity data. This level of customization ensures that the advice is relevant and practical, helping users achieve their fitness goals more efficiently.
- **User-Friendly Interface:** The web application is designed to be intuitive and accessible, allowing users to upload their data and interact with the system quickly. This accessibility encourages regular use and engagement, fostering a habit of continuous fitness tracking and improvement.

## Limitations

Despite the significant achievements, there are several limitations to address:

- **Local Backend Hosting:** Currently, the backend processes, including machine learning analysis and data processing, are hosted on a local machine. This setup limits the

scalability and accessibility of the application. Transitioning to a cloud-based infrastructure will enhance the application's performance and availability.

- **Broad Recommendations:** The LLM currently provides broad recommendations based on user data. More fine-tuning is required to deliver highly personalized and precise fitness advice. Advanced training techniques and larger datasets will help refine the model's recommendations.
- **Data Quality:** Ensuring high-quality data collection remains a challenge. Improvements in data validation and preprocessing techniques are necessary to enhance the accuracy of the analysis. Robust data collection protocols will ensure consistent and reliable input data.

## Future Work

Future developments in the FitLifeBot project will focus on several key areas:

- **Scalability and Accessibility:** The application's scalability and accessibility will be enhanced by transitioning the backend to a cloud-based infrastructure, enabling more users to take advantage of its features. Cloud hosting will provide better performance, reliability, and security.
- **LLM Fine-Tuning:** Further fine-tuning of the LLM will enable more specific and actionable fitness recommendations. Enhancing the LLM with advanced training techniques and more significant, diverse datasets will improve its ability to provide precise and personalized advice.
- **Enhanced Data Collection:** Implementing more robust data validation and preprocessing techniques will ensure higher quality and reliability of the input data. This will improve the overall efficiency of the suggestions.
- **Iterative Development:** Collecting user feedback will guide iterative improvements to the web application, ensuring it meets user needs and preferences. Regular updates and feature enhancements will be driven by user insights, making the platform more user-centric.
- **Integration with More Data Sources:** Expanding data integration to include other fitness tracking platforms and health data sources will provide a more comprehensive view of user health and activity. This integration will allow for a richer dataset and more holistic recommendations.
- **User Feedback Loop:** Creating an ongoing feedback loop with users will enable us to improve the application by considering real-world usage and user feedback. This method will guarantee that the application adapts to meet the evolving needs of its users.

## **Conclusion**

The FitLifeBot project successfully demonstrates the potential of leveraging Google Fit data and machine learning models to deliver personalized fitness recommendations. FitLifeBot offers users a comprehensive and interactive fitness experience by integrating detailed activity analysis with LLM technology. Addressing current limitations and focusing on future improvements will enhance the application's effectiveness and user engagement, making it a valuable tool for promoting healthier lifestyles. The project showcases a significant advancement in integrating fitness tracking and AI technologies, empowering users to manage their fitness progress and reach their wellness targets.

## References

References on human activity recognition and Google Fit data processing with proper working hyperlinks:

- Bulling, A., Blanke, U., & Schiele, B. (2014). A tutorial on human activity recognition using body-worn inertial sensors. *ACM Computing Surveys*, 46(3), 1-33.
  - <https://dl.acm.org/doi/10.1145/2499621>
- Lara, O. D., & Labrador, M. A. (2013). A survey on human activity recognition using wearable sensors. *IEEE Communications Surveys & Tutorials*, 15(3), 1192-1209.
  - <https://ieeexplore.ieee.org/document/6365160>
- Wang, A., Chen, G., Yang, J., Zhao, S., & Chang, C. Y. (2016). A comparative study on human activity recognition using inertial sensors in a smartphone. *IEEE Sensors Journal*, 16(11), 4566-4578.
  - <https://ieeexplore.ieee.org/document/7434149>
- Hossain, H. M. S., Khan, M. A. A. H., & Roy, N. (2017). Active learning enabled activity recognition. *Pervasive and Mobile Computing*, 38, 312-330.
  - <https://www.sciencedirect.com/science/article/pii/S1574119216301067>
- Chen, Y., Shen, C., & Wei, X. (2016). Attributed graph models for sensor-based human activity recognition. *Pattern Recognition*, 55, 87-99.
  - <https://www.sciencedirect.com/science/article/pii/S0031320316000364>
- Shoaib, M., Bosch, S., Incel, O. D., Scholten, H., & Havinga, P. J. (2015). A survey of online activity recognition using mobile phones. *Sensors*, 15(1), 2059-2085.
  - <https://www.mdpi.com/1424-8220/15/1/2059>
- Hammerla, N. Y., Halloran, S., & Plötz, T. (2016). Deep, convolutional, and recurrent models for human activity recognition using wearables. *arXiv preprint arXiv:1604.08880*.  
<https://arxiv.org/abs/1604.08880>
- Ordóñez, F. J., & Roggen, D. (2016). Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition. *Sensors*, 16(1), 115.  
<https://www.mdpi.com/1424-8220/16/1/115>
- Ravi, D., Wong, C., Lo, B., & Yang, G. Z. (2017). Deep learning for human activity recognition: A resource efficient implementation on low-power devices. In *2017 IEEE 14th International Conference on Wearable and Implantable Body Sensor Networks (BSN)* (pp. 71-74). IEEE.  
<https://ieeexplore.ieee.org/document/7936010>
- Banos, O., Galvez, J. M., Damas, M., Pomares, H., & Rojas, I. (2014). Window size impact in human activity recognition. *Sensors*, 14(4), 6474-6499.  
<https://www.mdpi.com/1424-8220/14/4/6474>