



**UNIVERSITA' POLITECNICA DELLE MARCHE**

**FACOLTA' DI INGEGNERIA**

---

Master's Degree in Biomedical Engineering

# **Monitoring of health and wellness status of elders users by telemedicine approach**

Supervisor:

Prof. **Lorenzo Scalise**

Candidate:

**Niccolò Sbaffi**

Co-supervisor:

Prof. **Sara Casaccia**

Dott. Ing. **Nicole Morresi**

Academic Year 2020/2021

Università Politecnica delle Marche  
*Department of Mechanical Engineering and Mathematical Sciences*  
Via Breccie Bianche — 60131 - Ancona, Italy



## Abstract

In recent years, the average life expectancy is considerably increased, creating a rapid growth of the aging population. Despite the medical advancement, the majority of the world's population is still affected by age-related health problems. Many research groups are studying pervasive solutions to provide independence and help to elderlies directly at their home. These solutions are based on the information acquired within a domestic environment using wearable and smart home sensors combined with AI based algorithms. In this context, the aim of this thesis is to analyse the data collected by commercially available smartwatch, worn by an elderly, and by a domestic sensor network, installed in single or multi resident homes, to monitor and improve the well-being of older adults by providing coaching solutions, predicting the wellness condition, or by identifying potential unusual behaviours. Three separate analyses are carried out on the available dataset to satisfy the proposed objectives. Firstly, the statistical analysis on smartwatch's data on a 2-weeks sliding period, help to classify the next available day as 'Normal' or 'Abnormal' and to provide coaching solutions. The second part is focused on the supervised machine learning analysis of the previous data such to predict the physical (PH) and mental (Mind) indices obtained through a daily survey. After the training, the Support Vector Machine with Radial Basis Function kernel provide the best accuracy of 100% for both Mind and PH indices. In the last part, unsupervised and supervised machine learning algorithms are applied to the unlabelled data of motion and light sensors of different homes to identify which sensor is the most informative to distinguish between 'abnormal' and 'normal' days at different time slot of the day.

# Contents

<b>CHAPTER 1: INTRODUCTION .....</b>	<b>1</b>
1.1 BACKGROUND AND MOTIVATION .....	1
1.2 THESIS OBJECTIVES .....	6
<b>CHAPTER 2: STATE OF THE ART.....</b>	<b>8</b>
2.1 WEARABLE TECHNOLOGY .....	8
2.2 SMART-HOME TECHNOLOGY FOR ELDERLY PEOPLE .....	12
<b>CHAPTER 3: STATISTICAL ANALYSIS OF ELDERLIES' SMARTWATCH DATA .....</b>	<b>14</b>
3.1 MEASUREMENT SET-UP .....	15
3.2 PARTICIPANTS.....	16
3.3 SMARTWATCH DATA .....	17
3.4 DATA ANALYSIS .....	17
3.4.1 <i>Pre-processing</i> .....	18
3.4.2 <i>Data processing</i> .....	19
3.4.3 <i>Comparison and classification</i> .....	19
3.4.4 <i>Coaching solutions</i> .....	21
3.4.5 <i>Impact of the smartwatch's uncertainty on the coaching solutions</i> .....	24
3.5 RESULTS .....	24
<b>CHAPTER 4: SUPERVISED MACHINE LEARNING ANALYSIS OF THE SMARTWATCH DATA.....</b>	<b>30</b>
4.1 METHODOLOGY .....	30
4.2 DAILY QUESTIONNAIRE.....	30
4.3 DATA ANALYSIS .....	32
4.3.1 <i>Exploratory Data Analysis</i> .....	32
4.3.2 <i>Questionnaire processing</i> .....	32
4.3.3 <i>Supervised Machine Learning</i> .....	33
4.3.3.1 Multi-class classification.....	35
4.3.3.2 Binary classification .....	36
4.3.3.3 Support Vector Machine .....	37
4.3.3.4 Decision Tree.....	39
4.3.3.5 Random Forest .....	41
4.3.3.6 MultiLayer Perceptron .....	42
4.3.4 <i>Performance metrics</i> .....	43
4.4 RESULTS .....	46

<b>CHAPTER 5: UNSUPERVISED MACHINE LEARNING ANALYSIS OF ELDERLIES' BEHAVIORAL DATA.....</b>	<b>50</b>
5.1 DOMOTIC DATA COLLECTION .....	50
5.2 EXPLORATORY DATA ANALYSIS.....	51
5.3 DATA PROCESSING.....	52
5.4 UNSUPERVISED MACHINE LEARNING.....	52
5.4.1 <i>Clustering Analysis</i> .....	53
5.4.1.1 K-Means .....	54
5.4.1.2 Principal Components Analysis .....	55
5.4.2 <i>Silhouette Score</i> .....	57
5.5 CLUSTERING INTERPRETATION USING SUPERVISED MACHINE LEARNING .....	57
5.6 RESULTS .....	58
<b>CHAPTER 6: CONCLUSIONS .....</b>	<b>60</b>
<b>BIBLIOGRAPHY .....</b>	<b>V</b>

## List of Figures

Figure 1: general architecture of a healthcare IoT system .....	3
Figure 2: example of wearable sensors worn on body parts.....	8
Figure 3: main functionalities of the iHealth Wave activity tracker. ....	16
Figure 4: plot of the Nr. of Steps data of the User 1. Missing or removed date are highlighted by red squares.....	19
Figure 5: plot of Nr. of Steps of the User 1 with the moving $\mu$ and $\sigma$ .....	20
Figure 6: example of coaching solution based on the Step count.....	22
Figure 7: example of coaching solution based on the Sleep Efficiency (%) data.....	23
Figure 8: example of uncertain days highlighted by 3 coloured bars. ....	24
Figure 9: plot of Nr. of Steps of User 1. ....	25
Figure 10: plot of the Sleep Efficiency (%) of User 1.....	25
Figure 11: plot of the Time in Bed (min.) of User 2. ....	26
Figure 12: plot of the Distance Travelled (km) of User 2.....	26
Figure 13: heat-map of Step dataset of User 1. ....	27
Figure 14: heat-map of Sleep dataset of User 1. ....	27
Figure 15: heat-map of Step dataset of User 2. ....	28
Figure 16: plot of uncertain days of User 1. ....	29
Figure 17: plot of uncertain days of User 2. ....	29
Figure 18: basic idea of SVM classifier.....	38
Figure 19: the training points mapped to a 3-D space where a separating hyperplane can be found. ....	38
Figure 20: Decision Tree scheme.....	40
Figure 21: Random Forest Scheme.....	42
Figure 22: MultiLayer Perceptron scheme. ....	43
Figure 23: binary Confusion Matrix.....	44
Figure 24: example of 3-class Confusion Matrix. ....	45
Figure 25: Accuracy score for the 3-class prediction of the PH index for the User 1.....	47
Figure 26: Accuracy score for the 3-class prediction of the Mind index for the User 1. ....	48
Figure 27: performance metrics for the 3-class prediction of the PH index for User 1. ....	48
Figure 28: performance metrics for the 3-class prediction of the Mind index for User 1. ....	49
Figure 29: schematic plan of one of the apartments with installed domotic sensors.....	51
Figure 30: K-Means algorithm. ....	55
Figure 31: transformation of high dimensional data to low dimensional data via PCA. ....	56

## List of Tables

<i>Table 1: characteristics of the participants.</i>	16
<i>Table 2: example of Step Data of the User 1.</i>	17
<i>Table 3: example of Sleep Data of the User 1.</i>	17
<i>Table 4: examples of missing or removed date in the Step data highlighted by red squares.</i>	18
<i>Table 5: results of the Statistical Analysis of User 1.</i>	28
<i>Table 6: results of the Statistical Analysis of User 2.</i>	28
<i>Table 7: percentage of uncertain days for User 1 and User 2.</i>	29
<i>Table 8: list of questions of the daily survey.</i>	31
<i>Table 9: example of dataset used for the supervised machine learning analysis.</i>	35
<i>Table 10: Accuracy score for the prediction of the PH index for the User 1.</i>	47
<i>Table 11: Accuracy score for the prediction of the Mind index for the User 1.</i>	48
<i>Table 12: Precision metrics for the prediction of the PH and Mind index for User 1.</i>	49
<i>Table 13: Recall metrics for the prediction of the PH and Mind index for User 1.</i>	49
<i>Table 14: F1 Score metrics for the prediction of the PH and Mind index for User 1.</i>	49
<i>Table 15: example of the Domotic Data of the Home nr. 3.</i>	51
<i>Table 16: example of Table for Time Slot 12 – 04 am of the Home nr. 3.</i>	52
<i>Table 17: Silhouette Score for the dataset.</i>	58
<i>Table 18: clustering rules for the House 4.</i>	59
<i>Table 19: clustering rules for the House 3.</i>	59
<i>Table 20: clustering rules for the House 5.</i>	59
<i>Table 21: clustering rules for the House 6.</i>	59
<i>Table 22: clustering rules for the House 7.</i>	60
<i>Table 23: clustering rules for the House 8.</i>	60
<i>Table 24: clustering rules for the House 10.</i>	60

# Chapter 1: Introduction

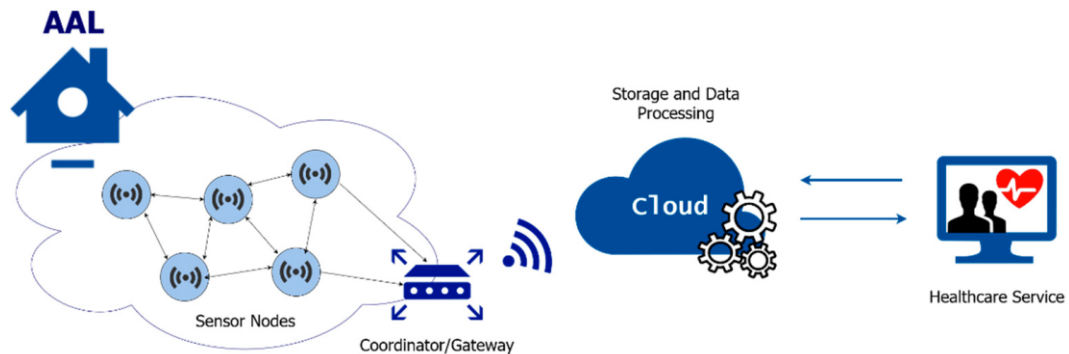
## 1.1 Background and motivation

Thanks to the recent advancements in medical science, diagnostic, and medicine, people's life expectancy is considerably increased, allowing people to live longer and healthier respect to previous generations, (1). The increase in life expectancy and the decrease in birth rate are predicted to generate soon a substantial elderly population. Worldwide, the number of people aged 65 years old or over was 727 million in 2020. In the next three decades, this number is estimated to more than double, achieving over 1.5 billion in 2050, (2). In terms of percentages, it is going to increase from 9.3% in 2020 to 16% of world's population by 2050, (3). Despite the development of medicine and the success of reducing mortality, a large part of the world's population is still affected by age-related health problems such as dementia, Alzheimer's, cardiovascular diseases, diabetes, chronic age-related diseases, limitations in physical activities and so on, (4). Therefore, the average life expectancy is increasing in parallel with the demand and costs for the healthcare services. This change in the demographic age represents an important human, social, and economic problem for the near future, creating many challenges for the world society and health-care system, since it is predicted that there will be a rise in the numbers of age-related diseases, followed by a rise in the costs and demand for healthcare services and a shortage of professionals trained to work with the elderlies, with more familiars taking the role of informal caregivers, and a rise in individuals unable to live independently, (5). Many negative consequences will affect the independent living and life quality of the elderlies while the sustainability of healthcare systems, given the rising cost of medical services, is going to impose a significant burden on the socio-economic structures of most countries, (6). In addition, many elderly people need routine assistance for their daily living and healthcare, which are mostly supported by the family members, (7). Formal paid care services offered by caregivers or elderly care institute have a high cost and thus are not affordable by a large part of the older population. Indeed, increasing the number of care providers to deal with the growing elderly population is not a feasible solution. On the contrary, "Telemedicine and Telemonitoring of elderly people is an actual challenge, investigated



to solve some of the problems linked to the constant growing of the average life expectancy”, (8). In the recent years there has been a growing interest into rapidly find and develop alternative solutions that will supply independence, assistance, and monitoring services for dependent people, like the elderlies, [ (9), (10)]. The development of these assisting systems is made possible by the convergence of advanced technologies in Machine Learning (ML) and pervasive computing. These innovative solutions are essential because the progressive deterioration in physical and cognitive abilities together with the aging diseases prevent the elderlies to live independently and well in their home and to perform their basic Activities of Daily Living (ADLs), [ (6), (11)]. Many research groups are studying proactive and pervasive technological solutions to provide help directly at the elderlies’ homes, by combining home environment and healthcare systems with the Internet of Things (IoT), (12). To gain knowledge about human health and well-being, the capacity to quantify and track the daily activities of individuals in their private home with sensors has become an important aspect, (13). Indeed, monitoring the well-being of individuals inside their home may help to improve their quality of life by promoting their independence and, what is more important, it will make possible the identification of any cognitive or physical decrease, that could lead to potential arising ageing diseases, [ (14), (15), (16)]. Smart technologies applied to pro-active health care, independent living, and active aging at home, introduced as ‘Ambient Assisted Living’ (AAL) technologies, can be a possible solution to the economic and societal challenges provoked by an aging population. These technologies have received increased attention from government, industry, and research institutes in recent years, (17). The advantages of AAL include saving long-term care costs while improving the quality of care, reduce family members and caregivers’ distress, and increasing the independence and overall quality of life of elderlies. In this way, medical facilities could be needed only for pathologies and emergencies. AAL is the inclusive term involving various ICT-based technologies based on the principles of ambient intelligence (Aml) such to supply age-friendly supportive and assistive environments to older adults and their caregivers, (18). Aml brings intelligence into everyday home by employing smart sensor networks, pervasive computing, and Artificial Intelligence (AI), making the environment sensitive, adaptive,

and responsive to the needs of any individuals living within it, (19). Sensors can be installed to provide a continuous information to caregivers without limiting the inhabitants' ADL, increasing their well-being and safety. Age-friendly environments will allow older people to age in a better physical and mental status, promoting their social inclusion and participation and helping them to maintain their autonomy and a good life quality.



*Figure 1: general architecture of a healthcare IoT system*

The classic ideas of Aml are applied to a new generation of assistive technologies for older adults, which are non-invasive and integrated into the environment, recognize the user and the situational context, adaptive (responsive to the user through learning), and anticipatory (anticipating the user's needs), (20). The introduction of smart sensors within a home environment capable of measuring explanatory quantities of the status of a subject together with algorithms capable of transforming data into useful information could allow to face the older population problem. AAL groups together many advanced technologies but with a particular attention on smart home technology, mobile and wearable technology, and assistive robotics. These technologies work along with advanced computational techniques including activity recognition, behavioural pattern discovery, anomaly detection, context modelling, planning, scheduling, and location and identity identification, [ (6), (20)]. An AAL environment includes components that are interconnected and communicate with each other: the environment and the user information are collected by the embedded sensors (sensing), all information is aggregated and sent to the cloud to be analysed and interpreted by the computational techniques that will decide the appropriate action (reasoning), and finally various types of actuators, intelligent interfaces and assistive devices facilitate

action and interaction with the user (acting) (as showed in Figure 1). These smart home environments have acquired a fundamental role in helping elderly, dependent people, and in alleviating the burden of health care workers. They are pervasive healthcare systems for elderly and dependent persons providing individual healthcare and social services such as nursing, rehabilitation, and health assistance directly at home, (21). Smart Home aims at monitoring and assessing the person's health condition and behaviour in performing ADLs allowing the elderly to live more independently and therefore enhancing his or her quality of life, (22). The objective is to detect any deterioration regarding the person's health and prevent major complications. Moreover, the system aims to maintain the dependency level and avoid, as long as possible, the usage of healthcare institutions, like nursing homes and hospitals. The integration of sensors in smart environments measuring the status of a person and the adoption of algorithms turning the data into useful information would allow to improve the elderly's well-being at their home, (23). The crucial need for the development of these technologies is justified by the mentioned growing of average life expectancy. Inside the AAL framework it is possible to place also the RESILIEN-T project, in which is developed this research activity, (24). RESILIEN-T is evolved in the framework of the AAL programs, whose main task is to improve the lifestyle and the quality of life of people with dementia and their caregivers [ (25), (26)]. The aim of the RESILIEN-T project is to develop an innovative and modular ICT solution for self-management of cognitive impairment, to reinforce the self-monitoring ability of affected people, with the purpose of slowing down their cognitive and behavioural decline. The project suggests the combination of elements already available on the market (i.e., tablet developed for older adults, existing wearable devices, lifestyle monitoring system) and newly developed elements (cloud platform and app) to allow People with Cognitive Impairment (PwCI) to self-control their conditions. The system architecture is represented by a modular, integrated, and open platform offering different services. The Cloud platform and portal components are realized to collect and process data, while the use of Application Programming Interfaces (API) allow the integration of additional devices in the architecture. This will allow the interoperability between the platform and devices. A tablet for older people is used as a gateway and data collector. The data

measured through the apps, wearable devices, and home monitoring systems are stored in the Cloud using the provided open API. The older user can use the tablet, with a specific user-friendly interface, to communicate with informal and formal caregivers, to send photo and message, to find events in his/her district, etc. The system architecture is designed to be modular, with 3 possible versions available now. The **Basic version** is composed of the tablet with ad-hoc modules to deliver coaching services relating to feeding, physical activity, cognitive exercises, and social interaction. The modules interact with a remote cloud-based platform to collect all the data generated from the user's interaction with the tablet, but also acting as a repository for the contents delivered to the user through the app. The tablet includes remote management by the informal carer supporting the PwCI. The **Plus version** is characterized by the basic system integrated with a wearable sensor (i.e., smartwatch) to collect mobility, physical activity, and other physiological data from the person. Lastly, the **Home version** is composed of the basic version integrated with on the market systems for lifestyle monitoring of people living alone with cognitive impairment. This is realized to collect new data from monitoring additional variables including sleeping habits and indoor mobility. Data analysis strategies, like ML and AI, could be used to add values to the data and information coming from the lifestyle monitoring sensors. The innovative nature of the proposed solution is linked to the integration of different components already available as stand-alone products. As already mentioned, the Plus Version of the RESILIEN-T project joins the Basic version of the system with a wearable device, like a smartwatch. A smart wearable device can be defined as a user worn accessory, with integrated electronic and computing technologies, that captures or reports some form of data, (27). Wearable data can include information about physical activity, movement, heart rate, temperatures, and blood pressure, that can be later used for analysis. Currently, smartwatch, smartphone and smart clothing are the conventional products joining together wearable technologies with care functions. All of them have attractive benefits for improving the well-being and delivering health information about the subject wearing them, (28). Wearable sensors have advantages over ambient sensors since they are generally smaller, cheaper, and always with the wearer, at any location, without requiring additional sensors to be placed in every room of an environment. The

need for wearable sensors in the healthcare sector continues to increase for applications such as long-term monitoring of patients in their homes, [ (29), (30)]. The smartwatches have the potential to change the health care and improve the well-being by supporting and assessing health in everyday living. Nowadays, they are simply available and familiar to most, allowing near-real time continuous monitoring of physical and physiological aspects, and enabling the communication between patients, family members, and health care providers. Therefore, such monitoring devices become an essential element for the older population that, with the increasing average of life expectancy, may develop mental and physical illnesses that will affect their independence and health. The Home version of the RESILIEN-T project combines the Basic Version with one of the available systems for the lifestyle monitoring, [ (31), (32)]. These systems are equipped with a network of smart sensors and actuators that collect data continuously and provide contextual information about the environment and the resident's activities. The aim is to monitor and assess the person's health condition and their behaviour in performing ADLs, make the healthcare services more economically sustainable, allow elderlies to live more independently, and improve the quality of life within their space. The sensor network is typically composed by Passive InfraRed (PIR) motion sensors, door or light sensors that are used for monitoring ADLs of the inhabitants, (9). These chosen sensors can be installed to provide a continuous information for caregivers without limiting the patient's ADLs. Any decay in the person's health could be detected in order to avoid or prevent major complications, (33). Moreover, the system aims to avoid, as long as possible, the delays of recourse to healthcare institutions (e.g., nursing homes and hospitals) by maintaining a high dependency level of elderlies at their home. The success of proposed health monitoring systems is based on the capacity to identify the context of supervised persons via the smart sensors installed in the home.

## **1.2 Thesis objectives**

The aim of the thesis is to analyse data collected from a wearable device and from sensors installed within a smart home environment, to monitor and improve the health and well-being of older users. As mentioned before, this thesis is born inside the framework of the RESILIEN-T project. Data, collected from wearable and domotic

sensors, are analysed through traditional methods (i.e., Statistical Analysis) and innovative approaches (i.e., ML and AI), with the final aim to identify and propose innovative strategies that will be at the basis of new healthcare services which could improve the well-being and the life quality of aging people directly at their home. In particular, this work is divided into three separate analyses, whose common objective is to try to measure and monitor the well-being of elderly users within their homes. The first part, named “Statistical analysis of elderlies’ smartwatch data”, provides a methodology to supply coaching solutions for elderlies during the day, starting from data collected through a commercial smartwatch, worn for a period of about 4 months. These data are analysed through a statistical approach performed over a 2-week period with a 1-day sliding window approach to re-new the time period under observation. In the second part of the thesis, the same data acquired by the smartwatch are used to predict a physical and mental index using various supervised ML algorithms. The 2 indices are obtained from the answers to a daily questionnaire provided by the elderly and containing questions relative to his or her physical and mental wellness. The aim is to verify the possibility of predicting the physical and mental status of the elderlies by training supervised ML algorithms using the acquired smartwatch’s data as features and the daily surveys’ indices as labels. The last part of the thesis deals with the unsupervised ML analysis of the domotic data relative to the ADLs of the elderlies within a home equipped with smart home sensors. The dataset is extracted from a specific domestic wireless sensor network, composed of PIR and Light sensors, installed in the most important room of single and multi-resident private homes of aging people. The aim is to analyse the behavioural data using the K-Means clustering algorithm to profile the different behavioural patterns for detecting possible anomalous behaviours among the different days of the dataset. The research work is organized as follows: Chapter 2 contains the state of the art discussing wearables and smart homes technologies for monitoring and improving the well-being of elderlies; Chapter 3 will explain the statistical analysis performed on the smartwatch’s data and will report the relative results; Chapter 4 will introduce the supervised ML analysis performed on the smartwatch’s data and daily surveys and will report the results; Chapter 5 will focus on the unsupervised ML analysis performed on the elderlies’ smart home behavioural data

and will explain the relative results; Chapter 6 is the Conclusive chapter that will discuss all the results obtained and the potential future works related to the thesis.

## Chapter 2: State of the Art

### 2.1 Wearable technology

The field of smart wearable devices has developed rapidly in recent years thanks to the success of mobile medicine, the development of new technologies and smart sensing with the tendency to get faster and smaller at the same time, and the increased popularity of personalized health concepts. These intelligent devices not only assist people in achieving a healthier lifestyle but also provide a continuous stream of health data for diagnosis and treatment of disease by actively recording physiological parameters and monitoring the metabolic state, (34). The terms wearable devices include miniaturized and mobile electronic devices, or computers with wireless communications that are incorporated into gadgets, accessories, or clothes, that can be put on the human body, as showed in Figure 2, (35). These smart wearable devices, given their portability and advanced computational efficiency, started to be employed in everyday life, not only within the medical context but also in the domestic one. Indeed, smart homes with ambient and wearable medical sensors, actuators, advanced connectivity, and ICT allow for the continuous and remote monitoring of the health and well-being of the elderly at a low cost.

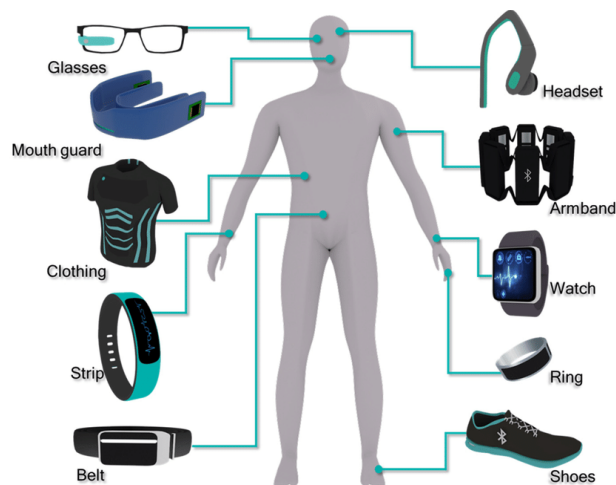


Figure 2: example of wearable sensors worn on body parts.

Older people are allowed to stay in comfortable homes while hospital caregivers can track the general health state of the elderly in real time with input and assistance from remote locations, (9). Internet of Wearable Things (IoWT) has emerged as a fundamental part of these smart environments. The IoT has turned basic wearable sensors into smart sensors, allowing objects to collect and exchange data through the internet with a manufacturer, operator, and/or other connected devices, without requiring human intervention, (36). Smartwatches, which are effectively considered wearable devices, have the potential to change the elderly health care by supporting and assessing their health in everyday life since they are easily accessible and familiar to most, allow near-real time continuous monitoring of physical and physiological quantities, and, most importantly, enable the communication between patients, family members, and health care providers. The adoption of wearable sensors for collecting qualitative and quantitative data, the capacity of measuring valuable quantities of the well-being of a subject, and the adoption of algorithms capable of transforming data into information could help to deal with the aging-in-place preferences of the elder's population. Sensors installed in the wearables can be divided into *inertial sensors* and *vital sign sensors* (or *biosensors*). The wearable inertial sensors can provide accurate features explicative of user's movement and body posture; *accelerometers* are the most frequently used sensors for activity monitoring. They can measure the value of acceleration along a sensitive axis and are particularly effective in monitoring activities including body motion such as walking, standing, sitting, or walking upstairs and downstairs. Due to their small size and low cost, accelerometers can be embedded into wrist bands, watches, and belts to monitor the user's activities, to detect possible falls, and wirelessly send data to mobile computing devices for further data analysis. The second category, which are the wearable *biosensors*, can collect many vital signs such as heart rate, blood pressure and skin temperature. These physiological parameters are fundamental for the monitoring of the elderly people's health condition. There are various biosensors used to measure the wide range of signals: *Electroencephalography* (EEG), *Electrooculography* (EOG), *Electromyography* (EMG), *Electrocardiography* (ECG), *pressure* sensors to monitor blood pressure, *CO<sub>2</sub>* sensors for the respiration, *thermal* sensors for monitoring body temperature. These physiological parameters help in the



monitoring of users' health status during the execution of ADLs. Based on the data collected from these biosensors, it is possible to provide some services such as early detection of disease, anomaly detection, and diagnosis decision-making. The information provided by these sensors are of undeniable utility, particularly in the case of elderlies. Most of the commercially available wearables allow to monitor many different quantities regarding older people's ADLs such as the number of steps per day, the sleep quality, or the calories burned during the activities performed. The importance of sleep and physical activity data, as representative of the actual and future health status, has been investigated by many research groups worldwide that have implemented different methodologies, all with the same research question, to see if this information could be used to improve the well-being and to predict the future wellness condition of the elderlies. The large amount of data obtainable by this technology is characterized by different determinants of physical, mental, and emotional states of the subject. With the help of the AI and ML in this context, it is possible to discover and explain the relationships existing among the quantities analysed. In (37), the authors proposed an integrated system to monitor the general wellness condition of elderly at the community level using two distinct electronic devices. The system framework included an electronic wearable wellness tracker and an all-in-one station-based health monitoring device that continuously monitored the elderly's overall activity and collected his or her daily vital signs, respectively. The designed personalized health monitoring scheme was developed for predicting the one-day-ahead wellness of elderly by choosing the most appropriate data mining models. The authors used 6 different classification methods: LASSO regression, Artificial Neural Network (ANN), Linear, Polynomial, Radial kernels Support Vector Machines (SVMs) and a Decision Tree (DT). The features of the models were a one-day lagged aggregated data from wearable wellness tracker and telehealth monitoring device, while the labels were a binary variable computed from a 10-point Health Index. DT method showed superior performance and it achieved a classification accuracy of 68.08%, with a recall of 64.02% and a precision of 81.65%, indicating its relative effectiveness in predicting elderly's health condition. In (38), the author proposed a one-day-forward forecasting method of wellness condition for community-dwelling elderly based on single lead short ECG

signals of 5 seconds that was used as input to the prediction model. A normalized Health Index was adopted as labels of the forecasting algorithm. The deep learning-based methods used were Long-Short Term Memory (LSTM) network and Bidirectional LSTM (BiLSTM) network. Meanwhile, the employed traditional ML-based methods included ANN and SVM. The BiLSTM achieved the best forecasting performance for the Health Index, whose recall, precision, accuracy, and F score were 92.51%, 91.48%, 93.21%, and 91.98%, respectively. In (39), the author evaluated the possibility of automatically predicting the well-being of elderly by gathering sleep data, activity, and time-away from home, using a medical sensor watch. In addition, the author collected both self-reported assessments of well-being from the participants and nurse evaluations on the participants' well-being. For predicting the well-being, different classifiers were trained and tested on different wellness labels. The classifiers used are K-Nearest Neighbour (KNN), Random Forest (RF), and SVM. Using the RF classifier, both the self-reported assessments and nurse evaluations could be forecasted with 94% accuracy and 90.64% accuracy, respectively. In (40), the authors presented an intelligent wearable system to monitor and predict mood states of elderly people during their daily life activities. The system was composed of a wristband device that records physiological and activity data of the users. In addition, the system used a mobile app for Ecological Momentary Assessments (EMA), whose input were used as ground truth (or labels) events for training the classifier. They built a mood classifier using the SVM with Radial Basis function. The best accuracies were: 90.05% for mood classification, 88.93% for happiness classification, and 87.21% for activeness classification. Based on the technologies investigated in the previous works, it is possible to project and develop digital 'coaching solutions' that are designed to improve the wellness of the elderly, by using smart sensors that are minimally invasive. Moreover, the possibility of predicting the elderly's well-being is an important research area, developed for providing proactive healthcare services to improve the life of aging people when a future deterioration is forecasted in advance. This thesis wants to propose solutions in both areas of research by combining a smartwatch, commercially available on the market, during everyday life with traditional statistical analysis and innovative ML algorithms to provide coaching solutions and to predict the physical and mental wellness of the user

wearing the device. The information used include the data monitored by most of the available activity tracker like daily Nr. of Steps, sleep efficiency, distance travelled, and calories Burned. Starting from these measurements, this work tries to validate the hypothesis that the physical and mental wellness of the elderly could be predicted from these readily available quantities.

## **2.2 Smart-Home technology for elderly people**

The idea of Smart Home arises with the introduction of network enabled devices and ultramodern electronic equipment usable in a domestic environment. The IoT and Smart Environments (SEs) have evolved from the simple homes into smart homes, changing the traditional approach of building devices, systems, services and transforming people's lifestyles. The rapid advances in the ICTs along with the increase of a large variety of affordable sensors and actuators have encouraged the growth of the IoT infrastructure. The development of SEs systems is one of the best results of IoT application, (41). These systems employ many small computational nodes to recognize and supply personalized services to the user that interact and exchange information with the environment, (42). The smart homes are the most typical examples of SEs and applications of IoT in the home environment, (43). The aim is to integrate different home-based objects with data communication capabilities to offer new or advanced functionalities to the inhabitants of the house. As stated in the introduction, due to the increase of the average life expectancy, alternatives to classical healthcare facilities should be found to efficiently take care of dependent people. Elderly people who live independently and have not developed severe diseases, need a solution that can monitor their ADLs in their home environment to ensure their safety, improve their lifestyle, and reduce the costs of health services and time. To meet these requirements, it is necessary to monitor these people using data extracted from specific sensor network. The use of the Smart Homes technologies can help older or disabled people, often with special needs but not only, to be autonomous, independent and continue to live in their own home for as long as possible. The ability of a smart home system to identify and predict the behaviour and wellness of its older occupants is strongly depending on the characteristics, number, and predisposition of the sensors involved,

as well as on the AI systems used to interpret the collected data. The combination of all these things provides the methodology for predicting or measuring the well-being among the elderlies, with the aim of intervene before any declines in physical and cognitive skills could affect them. In the recent years, many research studies were conducted to investigate the mentioned problem and most of these works focused on the possibility of improving the well-being of elderlies by monitoring their ADLs such to detect any possible abnormal behaviour that might require interventions. In (44), authors investigated the correlation between an older person's health status and his or her daily behaviour, using unsupervised ML algorithms, like the cluster analysis, to discover potential behaviour-related features from low-level sensors that could be easily installed in the home. The unlabelled movement sensors data were used to recognize different behavioural patterns of users belonging to 2 different clinical groups. The obtained results confirmed that the features considered in the study were able to group participants with similar behaviours from those with a different one. Furthermore, the proposed approach was able to identify participants with specific behavioural profiles that have the potential to inform more personalised caregiving and support. The methodology described could also be used to track participants' needs over time and fine-tune his or her care plan, optimising the care process. In (45), the authors investigated the utility of unsupervised ML and data visualisation for tracking changes in user activity, over time. This was done by analysing unlabelled data generated from passive and ambient smart home sensors, such as motion sensors. Results revealed that individual week-day data, considered over long periods, contain unique features that could be used to infer user activity levels and track any changes over the long term. The information discovered could be further utilised as part of a structured process or assessment protocol which could help to identify anomalies or changes in user activity. This could then be used for supporting carer-patient interactions, or even tracking the effectiveness of interventions and medication on the user's health condition as indicated by their activity or changes to routine over time. In (46), the authors designed an unobtrusive ADLs monitoring system to allow identification and prediction of elderlies' abnormal behaviour in the smart-home. The sensors used in this paper included movement and door entry point sensors. An individual inhabitant model was

learned from the collected data, which eventually represent the behavioural model of the user. When the user's pattern is learned, any anomalous behaviour was detected using clustering techniques. Moreover, the authors applied Recurrent Neural Networks (RNNs) to predict future sensor activity in order to warn a caregiver when an irregular behaviour would be expected in the near future. In (47), the authors used Self-Organizing Maps (SOMs) to identify irregularities related to abrupt changes in ADLs. The anomaly detection method was based on a composition of unsupervised classification technique, like the SOMs, and next activity prediction employing Markov model. A smart-home environment using the proposed method could adapt to behaviour of a user and could provide alarms to a carer or other responsible person if unusual activities were detected. Based on the reported works, a part of this thesis will be focused on the usage of unsupervised ML algorithms for analysing unlabelled sensor data acquired in a smart home equipped with a minimally invasive sensors network made by PIR and Light sensors. The aim of the analysis is to identify, at different time frame, which days will be considered different respect to the remaining ones. Moreover, the condition of diversity will be justified by the adoption of SML algorithm that will provide which sensor, in the different time slot, is the most informative one.

## **Chapter 3: Statistical analysis of elderlies' smartwatch data**

The main objective of the first part of the thesis is to profile the behaviour of an elderly subject by assigning a condition of 'Normality' or 'Abnormality' to his or her daily behaviour. This is made possible by statistically analysing the data collected through a wearable device. The results obtained can be used as input for a decision-making process that builds up a personalised Coaching Solution to help elderly people to extend and maintain their healthy quality of life via the promotion of good habits and behaviour. It is fundamental to encourage patient involvement, motivation, and empowerment in adapting to lifelong and healthy changes. An automatic, virtual coaching system could be ubiquitous, user-specific, constantly adapting, and collecting information about user's activity and habits while providing the necessary feedback in

the appropriate moment, (48). This thesis presents a proposed e-coaching system for promoting physical and mental well-being for elderly population by providing suggestions, congratulations, or alert messages to the users. The supplied messages will be chosen on the results of the statistical analysis of the smartwatch's data. The algorithms have been implemented in Python and the work took place in the workplace of Spyder.

### **3.1 Measurement Set-Up**

The wearable device used in the first part of this thesis is the activity tracker iHealth Wave by the company iHealth Lab Inc. (Europe), [ (49), (50)]. It is a connected fitness monitoring device that can measure the daily Number of Steps, the Distance covered (km), and the Calories Burned (kcal). Moreover, it can track the swim activities and parameters like the duration, the calories Burned (kcal) and it can recognise 3 types of swimming strokes. Another important functionality of the device is the ability to monitor the user's sleep quality by automatically recognizing the Sleep Start and Sleep End Time, measure the Sleep Efficiency (%), and report the Time spent in the Bed. So, it offers 24-hour personal monitoring system. All the data from the iHealth Wave can be synchronised via Bluetooth 4.0 with the iHealth MyVitals app. Figure 3 summarizes all its functionalities. The activity tracker is 30 m waterproof. It has an OLED touchscreen display, it weighs about 35 g, and it is powered by 3.7 V Li-ion 100mAh battery. With a full charge, the battery can last up to 7 days. The wearable tracker continuously perceives the movements of the body on a 3-axis accelerometer, so it is ideal for physical activities when the wrist on which the iHealth Wave is placed is an integral part of the movement performed. The data are recorded all the time it is worn and powered up, which allow the tracker to track the individual's movements in any direction. The processing occurs when the data is transferred to the software associated with the fitness tracker on the smartphone or laptop which it is synchronized with it. The app lets the individual know how many steps have been taken, the number of calories burnt, and the distance travelled. The calories burnt and the distance travelled are estimated from the daily number of steps using a mathematical formula and algorithm. The app allows the individual to interact with the information in a user-friendly manner. In addition to

the physical activity, the iHealth Wave analyses the sleep when the user wears it to go to bed. This analysis is done thanks to the 3-axis accelerometer and algorithm especially designed for the sleep analysis. Detected movements during sleep are analysed to provide different information available in a dedicated view of the app after the synchronization of the iHealth Wave: total sleep period, duration of light sleep, duration of deep sleep, duration of time spent awake, and Sleep efficiency. The iHealth Wave automatically enters in sleep mode, so the user has nothing else to do than to wear it on his or her wrist when he or she goes to bed. It also exits the sleep mode automatically when user wakes up and makes a few steps. The price of the smartwatch ranges between 40 and 70 €, making it an affordable device for most of the elderlies.



Figure 3: main functionalities of the iHealth Wave activity tracker.

### 3.2 Participants

One male and one female participant have been involved in the first part of the analysis. The 2 participants, being husband and wife, live in couple, with normal social status, have no pathologies and they can perform everyday activities normally. Both participants were provided with an iHealth Wave that was worn for a period of about 5 months. The following Table 1 summarizes the characteristics of the 2 subjects.

PARTECIPANT	SEX	AGE	CLINICAL DIAGNOSIS
Subject 1	F	63	No pathologies
Subject 2	M	70	No pathologies

Table 1: characteristics of the participants.

### 3.3 Smartwatch data

The smartwatch data used in the first part of the analysis are listed below:

- Daily Number of Steps
- Distance travelled (km)
- Calories burned (kcal)
- Sleep efficiency (%)
- Time in Bed period measured in hour and minutes (hh:mm)

From this point, the Nr. of Steps and the Distance travelled are together referred as Step data, while the Sleep Efficiency and the Time in Bed represent the Sleep Data. The data collection period, for this research, is from the 3<sup>rd</sup> of October 2020 to the 21<sup>st</sup> of February 2021. Data are daily collected and are downloaded from the Smartwatch App from which are extracted in .csv format. Table 2 and Table 3 reports the first few lines of the Step and Sleep data of the User 1, as an example.

Date	Nr. of Steps	Distance (Km)	Calories (Kcal)
2020-10-03	1378	1.04	1344
2020-10-04	1706	1.28	1354
...	...	...	...

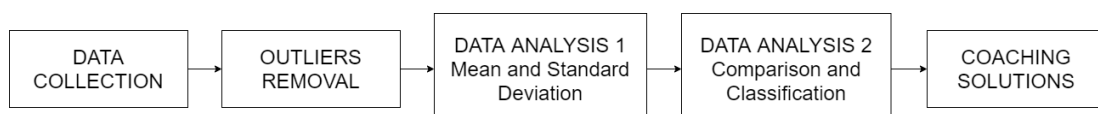
Table 2: example of Step Data of the User 1.

Date	Time in bed (h, min)	Sleep efficiency (%)
2020-10-03	11 hours 10 minutes	91
2020-10-04	8 hours 20 minutes	95
...	...	...

Table 3: example of Sleep Data of the User 1.

### 3.4 Data analysis

In the following paragraphs, a description of the methodology adopted to perform the Statistical Analysis of the smartwatch data is given. The data collection period includes most of the days between the 3<sup>rd</sup> of October 2020 and the 21<sup>st</sup> of February 2021. The methodology adopted can be divided in 5 steps as follows:





### 3.4.1 Pre-processing

After the data collection, the first step of the smartwatch's data statistical analysis is the removal of all the possible outliers, defined as samples that are exceptionally far from the mainstream of the data. The data collected are pre-processed to remove all the possible days that contain values identified as outliers. For what concerns the Step data, a daily Nr. of Steps equal to 0 is to be considered as an outlier, since users are expected to move or walk every day; in addition, a step count equal to zero can also mean that the smartwatch was not worn by the user, thus providing useless information. The same principle is applied to the distance travelled and to the calories burned. Table 4 reports some examples. Regarding the sleep efficiency and time in bed data, if one of the 2 values referred to the same day is equal to 0 or the time in bed value is over a threshold of 14 hours, that specific day needs to be classified as outlier. As consequence of the pre-processing step, all the days identified as outliers are deleted. As example, Figure 4 displays the entire period of data collection for the Nr. of Step features of the User 1. On the Y-axis is reported the corresponding component of the dataset, while each coordinate of the X-axis represents a day. The lack of solid line between two or more consecutive date means that these dates have been removed because are considered outlier. Moreover, some dates are missing because the smartwatch's malfunctioning provokes the non-acquisition of the data. The total number of days of User 1 after the first step is 133 and 130 for the Step and Sleep data, respectively.

Date	Nr. of Steps	Distance (km)	Calories (kcal)
2020-10-22	6449	4.9	1558
2020-10-23	3086	2.33	1413
2020-10-24	112	0.08	556
2020-10-26	3506	2.65	1422
...	...	...	...
2020-11-11	0	0	1297
...	...	...	...

Table 4: examples of missing or removed date in the Step data highlighted by red squares.

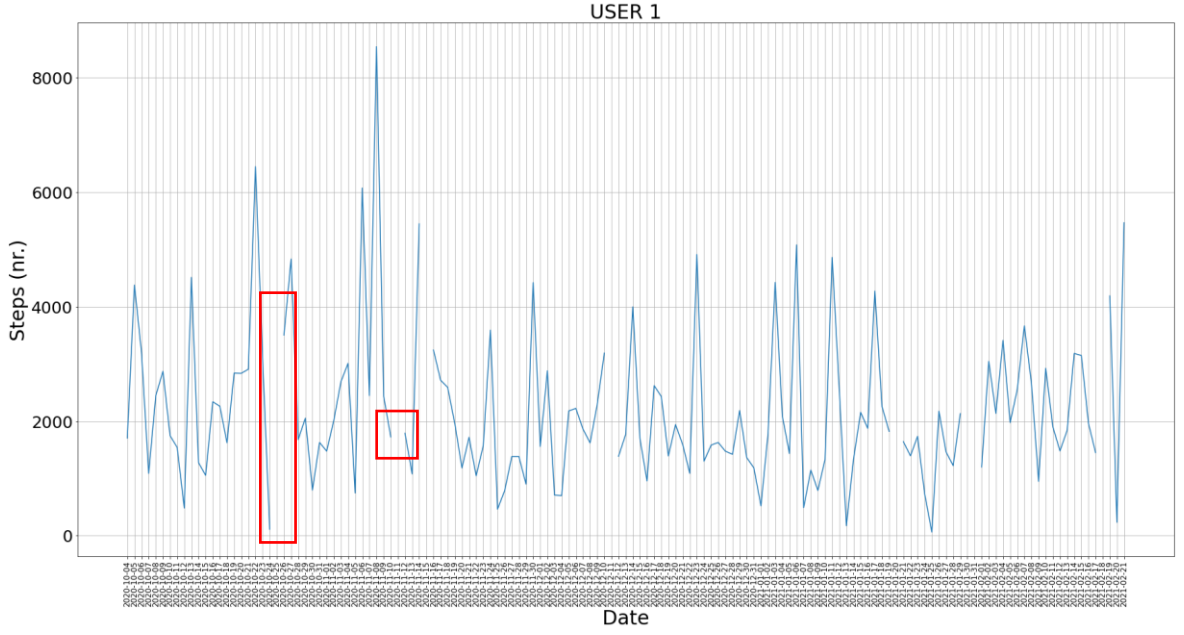


Figure 4: plot of the Nr. of Steps data of the User 1. Missing or removed date are highlighted by red squares.

### 3.4.2 Data processing

The second step of the statistical analysis is the computation of the Mean ( $\mu$ ) and Standard Deviation ( $\sigma$ ) values for each component of the Step and Sleep Dataset. The two variables are calculated over a 2-weeks period with a sliding window of 1 day. The mathematical equations of the  $\mu$  and  $\sigma$  are reported in the following lines:

$$\text{Mean } (\mu) = \frac{x_1 + x_2 + x_3 + \dots + x_N}{N} = \sum_{i=1}^N \frac{x_i}{N} \quad (1)$$

$$\text{Standard Deviation } (\sigma) = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{N}} \quad (2)$$

- $x_i$  are the data values/observations
- $N$  is the number of considered observations ( $N=14$ )

### 3.4.3 Comparison and classification

The third step of the statistical analysis is the comparison between the data collected by the smartwatch everyday with the range  $\mu \pm \sigma$  computed over the just preceding 2-weeks period. When the comparison is completed, the next step is the classification of each day in three different categories. The classification rules are based on the inclusion or non-inclusion of the smartwatch's data within the range  $\mu \pm \sigma$  estimated on the

previous 14 days. The following lines summarize the classification rules adopted for the assignment of each day to one of the 3 possible categories ‘Negative Abnormality’, ‘Normality’, ‘Positive Abnormality’:

- 1) The Rules nr. 1 assigns a condition of ‘Negative Abnormality’ to the day when the smartwatch’s data under analysis is lower than  $\mu - \sigma$  of the previous 2-week period.
- 2) The Rules nr. 2 assigns a condition of ‘Normality’ to the day when the smartwatch’s data under analysis is within the range  $\mu \pm \sigma$  of the previous 2-week period.
- 3) The Rules nr. 3 assigns a condition of ‘Positive Abnormality’ to the day when the smartwatch’s data under analysis is higher than  $\mu + \sigma$  of the previous 2-weeks period.

The procedure is repeated for each component of the data collected by the smartwatch iHealth Wave of User 1 and 2. The comparison and classification processes are repeated for each day, for this reason the values of  $\mu$  and  $\sigma$  are updated cyclically by selecting a 2-weeks period that is renewed with a sliding window of 1 sample. The sliding process is repeated until the last day of the dataset is compared and classified using the previous 2-weeks interval. Figure 5 display the just-mentioned approach applied to the Nr. of Step of the User 1. On the Y-axis is reported the Nr. of Step of each day, while each coordinate of the X-axis represents a day of the period analysed. The gap between the days of the X-axis indicates that the day is missing or is deleted because it is considered an outlier.

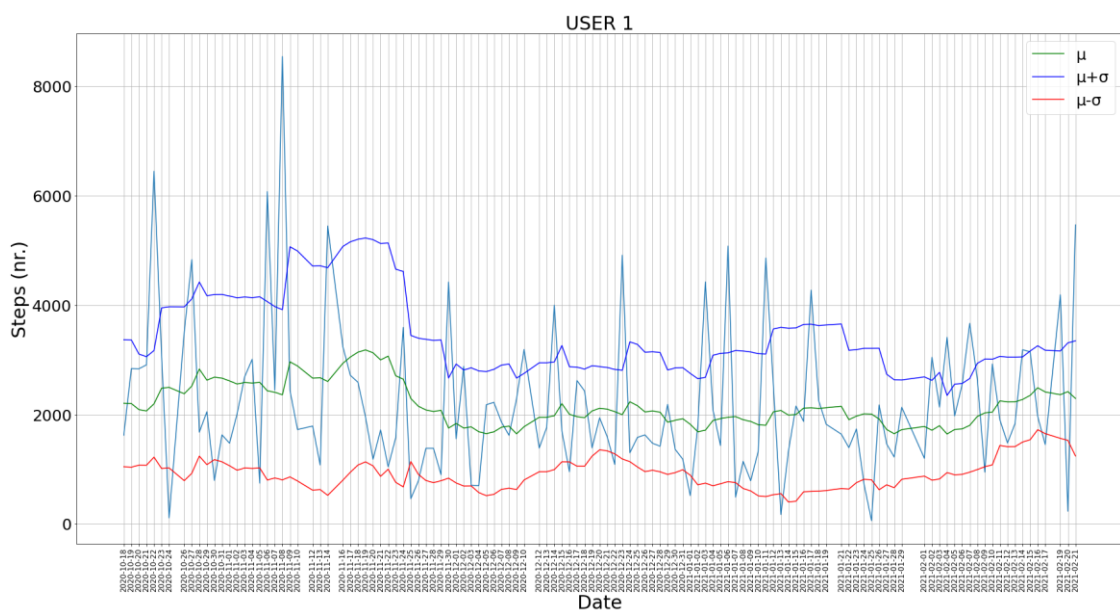


Figure 5: plot of Nr. of Steps of the User 1 with the moving  $\mu$  and  $\sigma$ .

### 3.4.4 Coaching solutions

The last step of the Statistical Analysis is the development of e-coaching solutions based on the results obtained in the Comparison and Classification part, previously described. An automatic, digital coaching system could be ubiquitous, personalized, constantly adapting and collecting information while providing the necessary feedback. Nowadays, technology-based interventions to improve healthcare and to promote wellness have become reality and they represent an important opportunity to promote healthy aging at a larger scale, (51). The advances and diffusion of mobile technologies (i.e., smartphones and IoT object) paved the way for the development of virtual coaches or e-coaches, which can support, complement, and possibly replace human coaches in health context. This thesis presents an e-coaching system for promoting physical activity and improve the sleep habits of the elderlies. Figure 6 and Figure 7 summarize the methodology adopted for the development of the e-coaching solutions. The 2 flowcharts describe the main passages of the algorithm's workflow that can implemented in the tablet or smartphone provided to the elderly. The first part half of the diagram is based on the statistical analysis applied to the smartwatch's data. The different boxes perform all the necessary steps to arrive at different outputs given by the flowchart. The classification results are used as decision criteria between 3 main types of messages that include suggestions in case of 'normality', congratulations in case of 'positive abnormality', or alert messages in case of 'negative abnormality'. The different texts are user specific, reporting the name of the elderlies wearing the smartwatch. Moreover, every time that the virtual coach recognises the 'negative abnormality' conditions of the previous day, it will ask to the elderly if there is a physical or mental problem affecting his or her wellness by limiting the ADLs or affecting the ability to sleep well on a regular basis. In case of positive answer, it will also ask if a medical examination could be necessary to solve the issue. The integration of these digital coaching within the smartphone or tablet makes possible the last functionality. On the contrary, in case of negative answer, the flowchart ends the process without any medical examination.

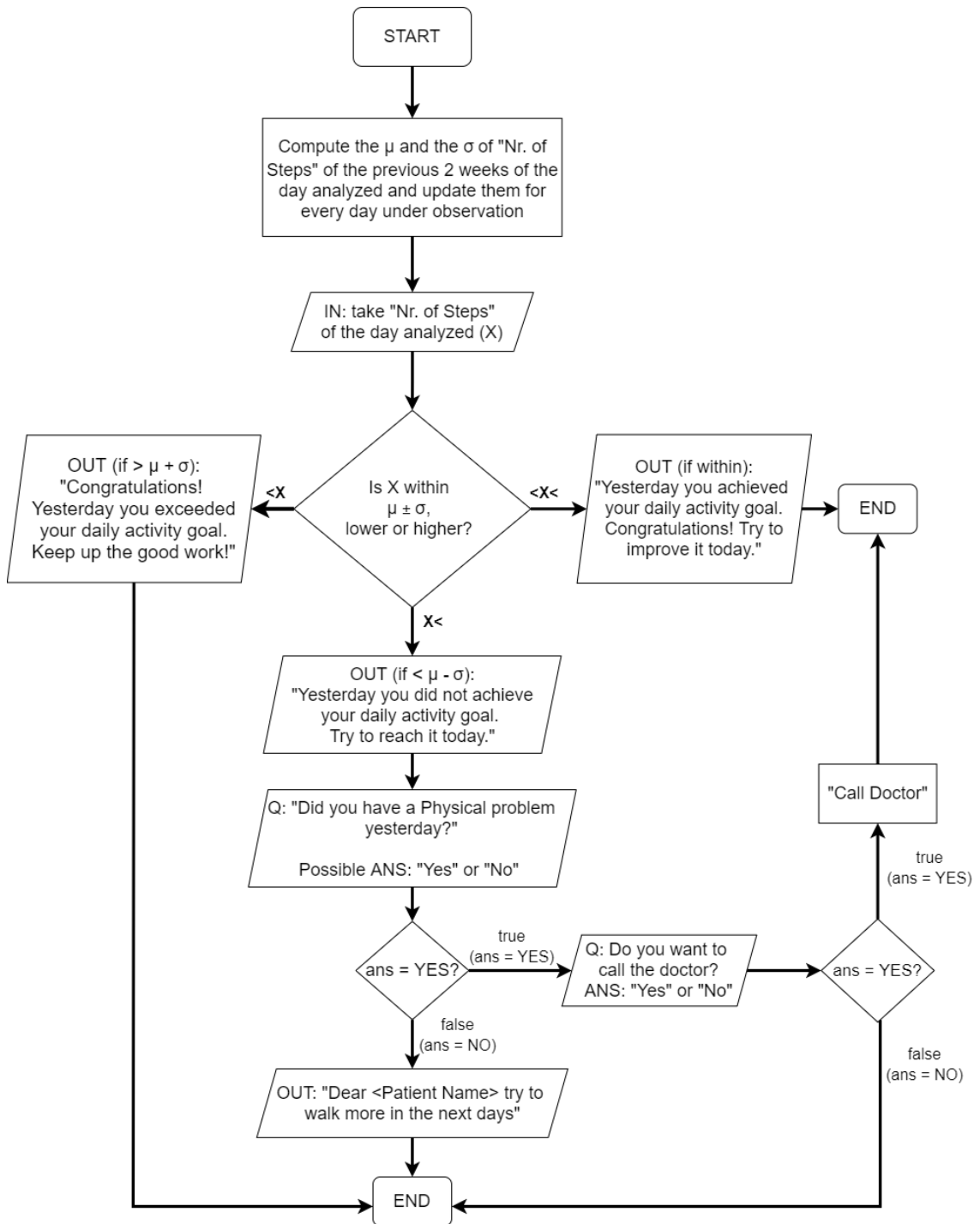


Figure 6: example of coaching solution based on the Step count.

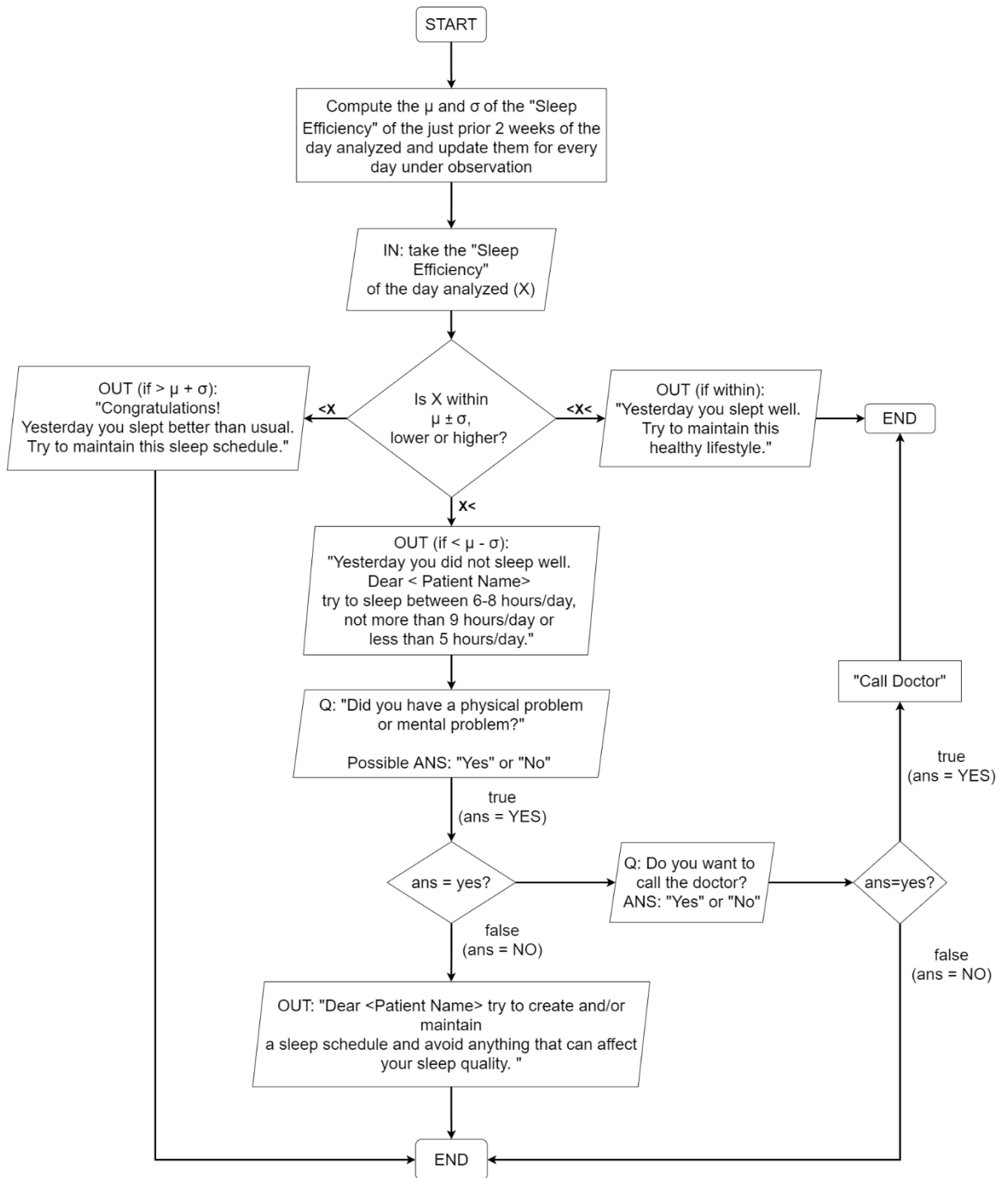


Figure 7: example of coaching solution based on the Sleep Efficiency (%) data.

### 3.4.5 Impact of the smartwatch's uncertainty on the coaching solutions

The iHealth Wave smartwatch has a measurement uncertainty on the count of the daily number of steps of  $\pm 4.7\%$ , for this reason this part of the statistical analysis is focused only on the signals of daily physical activity. This uncertainty should therefore be considered, because in the context of coaching, it can happen that this value can lead to misleading interpretation, providing an inadequate coaching solution, (52) . In fact, it can happen that the classification of the day under analysis is in two possible categories when the upper or lower limit of its range is higher than  $\mu + \sigma$  or lower than  $\mu - \sigma$  of the previous 2-weeks period. When this eventuality occurs, the coaching solutions provided to the elderlies may be incorrect and not truthful, bringing to a wrong message provided to the elderly by the e-coaching system. The methodology used for evaluating the impact of the smartwatch's uncertainty on the resulting coaching solution is performed as follows: first, every sample collected from the Nr. of Step is perturbed with the 4.7% uncertainty of its actual value; this perturbation will give rise to a confidence interval in which the measurement can be found; if the measured value and the result of the perturbation belongs to same category, the value is not considered as uncertainty in the classification process, otherwise, if the measured value, and the resulting perturbation are included in two different categories, the day is considered as 'uncertain' day. An example of the methodology is showed in Figure 8.

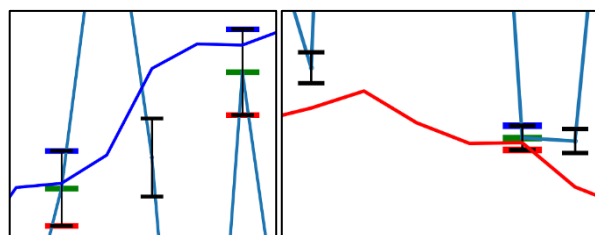


Figure 8: example of uncertain days highlighted by 3 coloured bars.

## 3.5 Results

The statistical analysis described in the previous chapter is applied to the data acquired by the iHealth Wave, worn by User 1 and User 2. Figure 9 shows the moving average applied to a period of 2-weeks with a 1-day sliding window approach computed for the

Nr. of Steps of User 1. The green line represents the value of the  $\mu$  over the previous 2-week period, while the blue and red ones are the  $\mu + \sigma$  and the  $\mu - \sigma$  over the same period, respectively. Same approach is applied to the remaining data collected via the smartwatch. As example, Figure 10 displays the same approach applied to the Sleep Efficiency of User 1. Figure 11 and Figure 12 display the methodology for other two features of User 2.

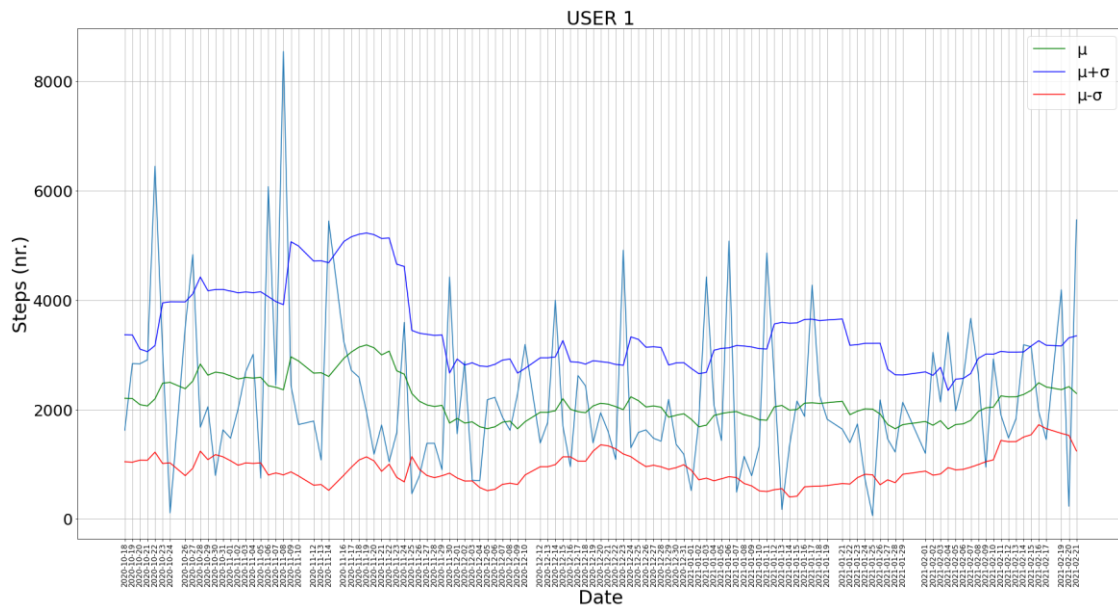


Figure 9: plot of Nr. of Steps of User 1.

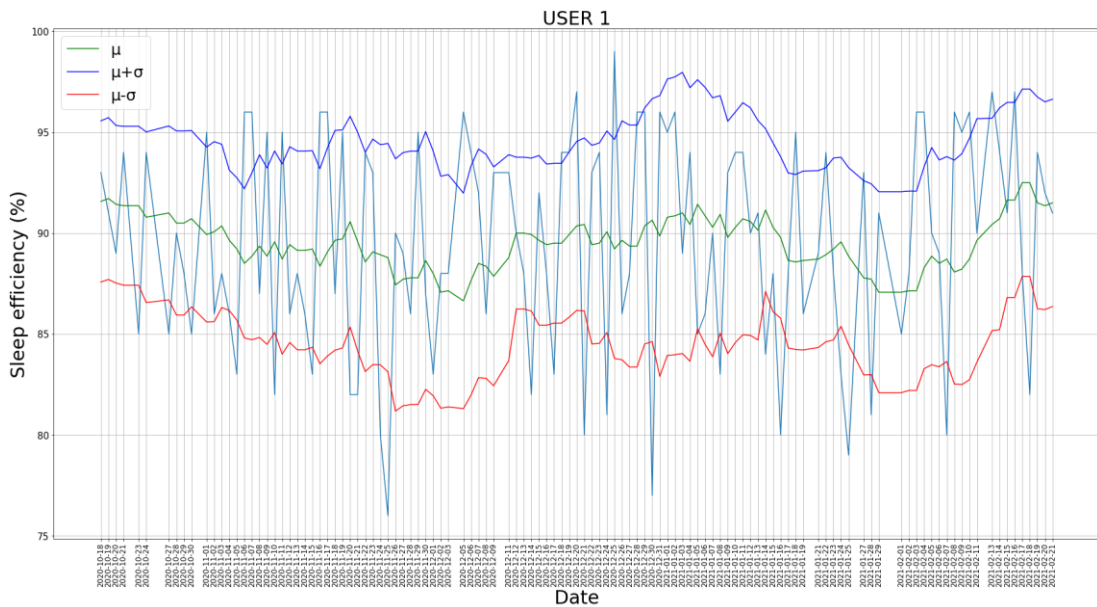


Figure 10: plot of the Sleep Efficiency (%) of User 1.



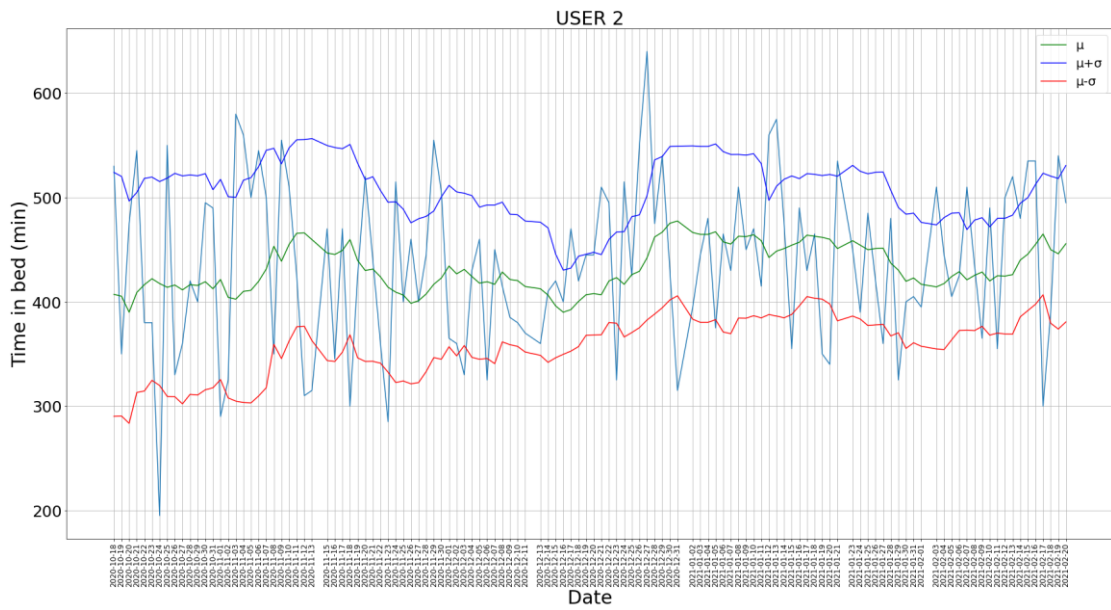


Figure 11: plot of the Time in Bed (min.) of User 2.

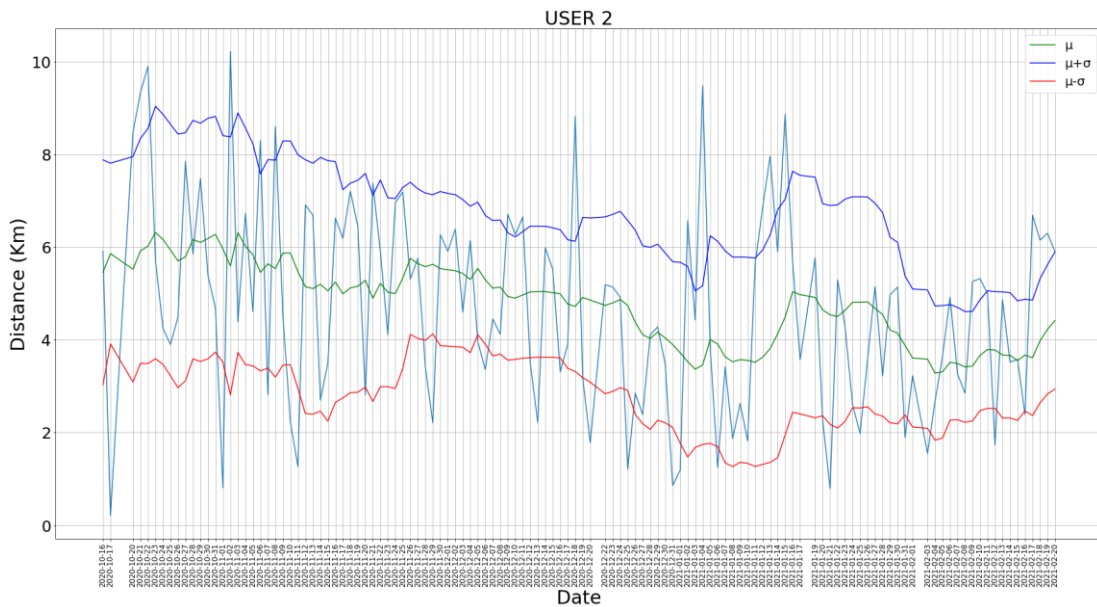


Figure 12: plot of the Distance Travelled (km) of User 2.

Figure 9 to 12 allow to visualize, for each parameter collected from the smartwatch, which days are going to be classified as ‘Normality’, ‘Positive Abnormality’, or ‘Negative Abnormality’. All the days within the range  $\mu \pm \sigma$  are classified within the category ‘Normality’, while the days outside this range are classified as ‘Positive Abnormality’ or ‘Negative Abnormality’, respectively. The following heat-maps display the classification results for the days of User 1.

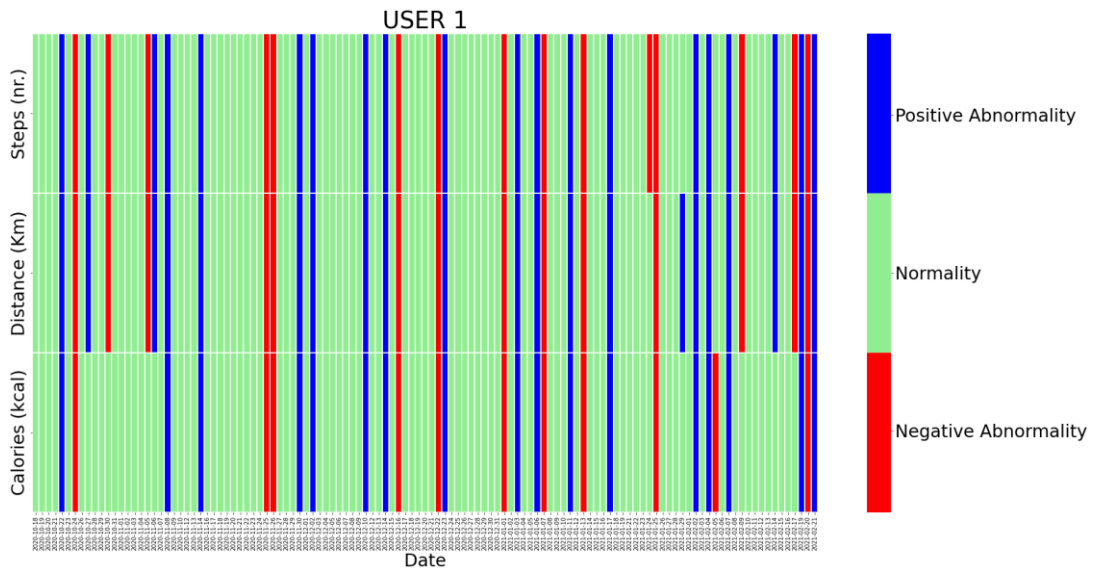


Figure 13: heat-map of Step dataset of User 1.

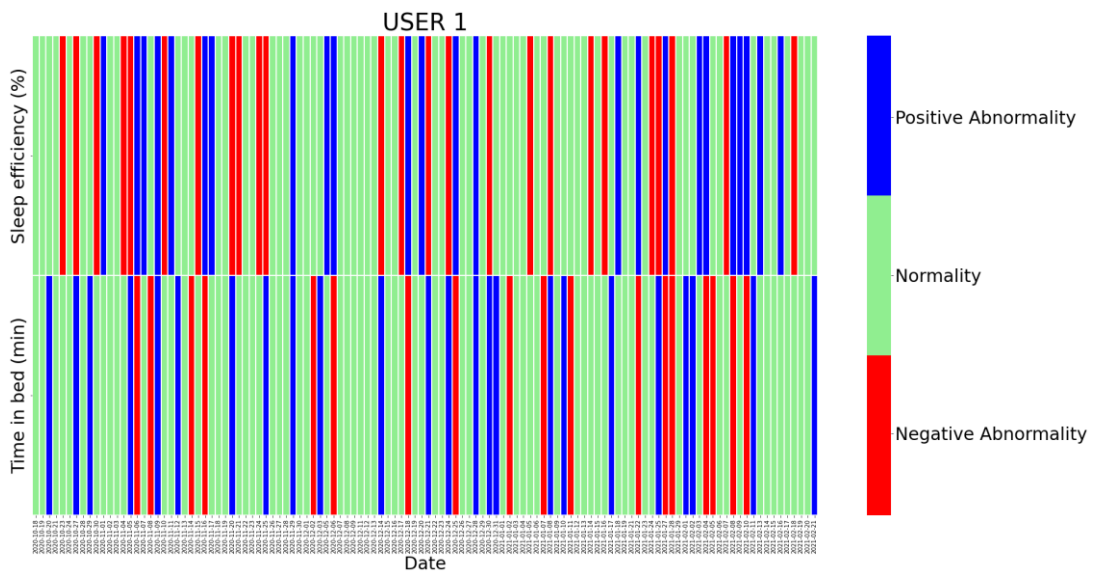


Figure 14: heat-map of Sleep dataset of User 1.

The heat-maps, obtained also for User 2, are displayed in the next figures. It is possible to note that the classification of the day based on the Nr. of Steps, Distance, and Calories is mostly the same for both Users. That is because from the steps count it is possible to compute, via mathematical equations, the calories burned, and the distance travelled. Therefore, most of the classes assigned to the various days correspond among them. The following tables report the percentages of days classified within the ‘Positive Abnormality’ and ‘Negative Abnormality’ classes out of the total dataset. Table 5 and 6 summarized the results are obtained from the statistical analysis of User 1 and User 2.

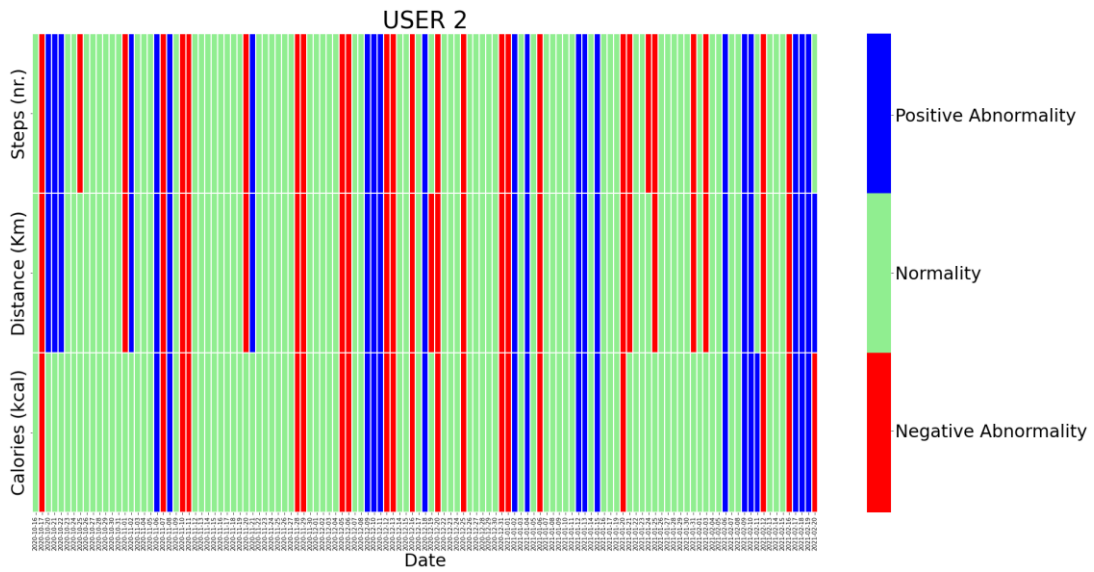


Figure 15: heat-map of Step dataset of User 2.

USER 1	% Negative Abnormality	% Positive Abnormality
Steps (nr.)	17	23
Distance (Km)	16	24
Calories (kcal)	13	20
Step dataset (tot.)	119 days	119 days
Sleep eff. (%)	29	27
Time in bed (min)	20	27
Sleep dataset (tot.)	116 days	116 days

Table 5: results of the Statistical Analysis of User 1.

USER 2	% Negative Abnormality	% Positive Abnormality
Steps (nr.)	33	27
Distance (Km)	31	28
Calories (kcal)	24	22
Step dataset (tot.)	123 days	123 days
Sleep eff. (%)	25	24
Time in bed (min)	24	35
Sleep dataset (tot.)	121 days	121 days

Table 6: results of the Statistical Analysis of User 2.

The classification results based on the steps count can be partially wrong and not reliable given that the smartwatch's measurement uncertainty on this data is about 4.7%. This means that the count of the Steps can be included in the range between its actual value  $\pm 4.7\%$  of the latter. This can provoke the classification of the day in 2 different categories when the actual value is near the separation boundary of the  $\mu + \sigma$  or  $\mu - \sigma$ . Figure 16 and 17 help to understand this eventuality. The whiskers lines above and below each steps count represent the possible range of values that this latter can assume. The days

belonging to 2 classes are highlighted by 3 coloured asterisks. Table 7 reports the number of uncertain days out of the total days for User 1 and User 2.

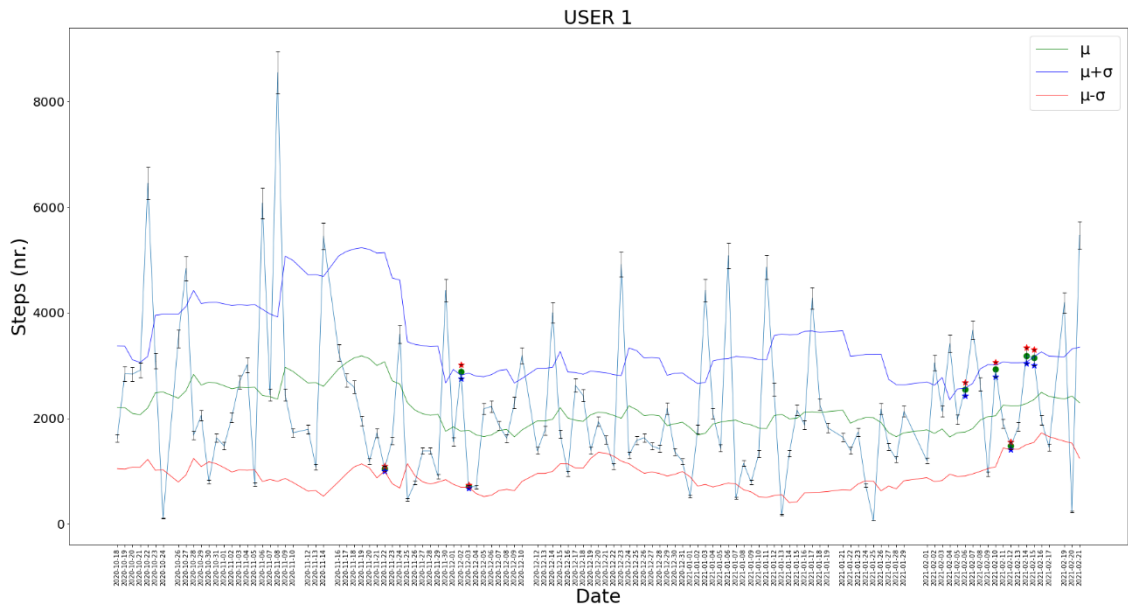


Figure 16: plot of uncertain days of User 1.

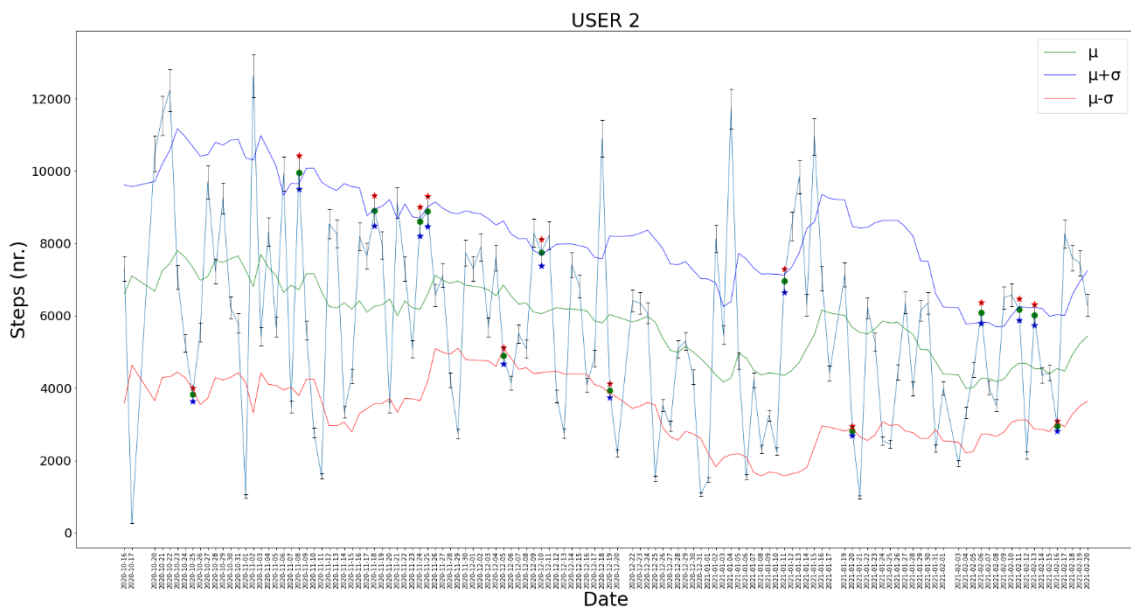


Figure 17: plot of uncertain days of User 2.

USER	Total days	Uncertain days	Error
User 1	119	8	6.7 %
User 2	123	14	11.4 %

Table 7: percentage of uncertain days for User 1 and User 2.

# Chapter 4: Supervised Machine Learning analysis of the smartwatch data

The main objective of the second part of the thesis is to propose a methodology that can predict the physical and mental well-being of an elderly subject using the data acquired by the iHealth Wave smartwatch and a daily self-evaluation survey about the health and mental status of the subject. The prediction of the user's well-being is possible by training Supervised ML (SML) algorithms using the questionnaire's answers as reference data. The idea behind this analysis is that the changes in the behavioural patterns measured with wearable sensors could be associated with the variation of the user's well-being.

## 4.1 Methodology

The activity tracker iHealth Wave (previously described in Paragraph 3.1) by the company iHealth Lab Inc. (Europe), [ (49), (50)] is used for acquiring the data employed in the SML analysis. The smartwatch data used in the second part of the analysis are listed below:

- Daily Number of Steps
- Distance travelled (km)
- Sleep efficiency (%)
- Time in Bed period (hh:mm)

Unlike the previous part, the daily Calories burned are not considered for the SML analysis since they are computed from the daily Nr. of Steps via a mathematical equation and so they do not provide additional information. Only the data provided by one female user, so User 1, is utilised for the purpose of this study.

## 4.2 Daily Questionnaire

Normally the well-being of users can be verified through a series of questionnaires, as reported in literature, [ (53), (54), (55)]. The user was asked to daily fill the questionnaire to report about her physical and mental status. From the survey it is possible to extract two numerical indices representing the self-evaluation of user's well-being. The daily

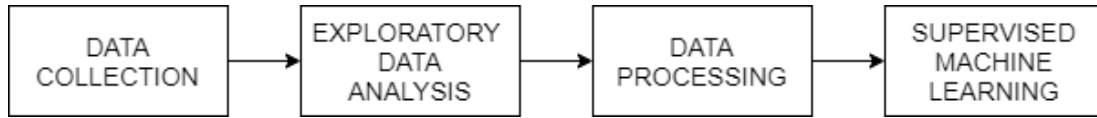
survey is created as a subset of the MOS Short-Form-36 (SF36) questionnaire, (53) . The number of questions is reduced to make the questionnaire easier and quicker to be completed, decrease the daily effort to fill it in and improve the acceptability. The survey (reported in Table 8) consisted of ten multiple-choice items asking to rank general health perception, functional status (i.e., housekeeping activity, physical activity, role limitation due to physical issues, general physical perception), mental wellness (i.e., mental health, role limitations due to mental problems) and guests' presence. For each question, older user can answer using an ordinal scale ranging from 1 to 3, where a rank of 3 meant that users were able to perform a great amount of activities during the day and no limitations due to mental or physical conditions occurred, therefore they had a positive self-perception, while a rank of 1 meant that they performed no activities during the day or that limitations due to mental or physical conditions prevented them from carrying out everyday tasks, which resulted in a negative self-perception. The last question has a binary answer: 'yes' or 'no'.

Nr.	Question
1	In general, how would you describe your health today?
2	Have you performed moderate activities (e.g., housecleaning, cooking, etc.) today?
3	Has your health status limited you in carrying out these moderate activities?
4	Have you carried out physical activities (e.g., walking, climbing stairs, etc.) today?
5	Has your health status limited you in carrying out these physical activities?
6	In general, how would you describe your mental health today?
7	Has your mental status affected your daily routine today?
8	Have you felt well physically (e.g., no aches, pains, etc.) today?
9	Have pains limited you in carrying out your daily activities today?
10	Have you received visits today?

*Table 8: list of questions of the daily survey.*

## 4.3 Data Analysis

The Methodology adopted in the second part of the thesis can be divided in 4 steps as follows:



The data collection period includes the days between the 3<sup>rd</sup> of October 2020 and the 21<sup>st</sup> of February 2021.

### 4.3.1 Exploratory data analysis

Like the previous statistical analysis, the first part of the supervised learning is the removal of all the possible outliers. The same principles, described in the Paragraph 3.4.1, are applied to all the data used for the SML analysis. ML algorithms are sensitive to the range and distribution of attribute values. Data outliers can ruin and mislead the training process resulting in longer training times, less accurate models, and ultimately poorer results.

### 4.3.2 Questionnaire processing

Questionnaire's answers were stored in an Excel file. Each row of the file reports the date of fulfilment, the user's name/surname, and his/her relative answers. Each answer is mapped to the range 1 to 3 so that the worst answer is represented by the value 1, the middle answer by the value 2, while the best answer is equivalent to a value of 3. The binary answer to the 10<sup>th</sup> question is converted to a value of 0 or 1 for the 'no' or 'yes' answer, respectively. When all the answers are mapped to a numerical value, the following step is the Normalization of their value using the following formula:

$$q_i = \frac{q_i - \min}{\max - \min} \quad (3)$$

where  $q_i$  is the  $i$ -th answer ( $i= 1, \dots, 9$ ),  $\max = 3$ , and  $\min = 1$ .

The questionnaire file of the user does not contain all the days of the analysis period since the user can decide to fill in the questionnaire every day if her physical and mental wellness change day by day or to skip some days between one survey and the following

one when the physical and mental status do not change for some days. In order to have a questionnaire file as much complete as possible it is necessary to add all the missing days between one questionnaire date and the next one. Each added day and the relative empty questionnaire are filled by copying the answers of the previously available day/questionnaire. Therefore, the result of this step is a user-specific file containing for each day/questionnaire her relative answers mapped to a numerical value normalized in the range from 0 to 1. Starting from the numerical mapping, it is possible to compute the Physical (PH), and Mind index for each day. The 2 indices are calculated according to the following formulas:

$$PH\ INDEX = \frac{quest\ 2 + quest\ 3 + quest\ 4 + quest\ 5}{4} \quad (4)$$

$$MIND\ INDEX = \frac{quest\ 6 + quest\ 7}{2} \quad (5)$$

where  $quest_i$  is the answer to the  $i$ -th questions, with  $i = 2, \dots, 7$ .

The range of numerical values attributable to the 2 indices is from 0 to 1. The user-specific Questionnaire, Step, and Sleep data are merged based on the common date available for all the 3 items such that the final data used for the analysis contains the Questionnaire, Step, and Sleep data for each day. When one of the 3 data is not present for a specific date, this day must be discarded. The total number of days after this step is 124.

### 4.3.3 Supervised Machine Learning

ML approaches are traditionally divided into three categories, depending on the nature of the "signal" or "feedback" available to the learning system:

1. Supervised learning: the computer is provided with example of inputs and their desired outputs, that can consist of numeric values or string labels, and the goal is to learn a general rule that maps inputs to outputs. So, the ML algorithm is trained on small training dataset of labelled data. It is very similar to the final dataset and supply the algorithm with the parameters required for the problem.
2. Unsupervised learning: unlabelled data are given to the learning algorithm, leaving it on its own to find hidden structure and patterns in its



input. Relationships between data points are gathered by the algorithm in an abstract manner, without any inputs from human beings. Unsupervised learning can be a goal in itself (discovering hidden patterns in data) or a means towards an end (feature learning).

3. Reinforcement learning: a computer program interacts with a dynamic environment in which it must perform a certain objective. As it navigates its problem space, the algorithm is provided feedback that is analogous to rewards, which it tries to maximize. So, it is characterized by an algorithm that improves itself and learns from new situations using a trial-and-error method. Favourable outputs are encouraged, or 'reinforced', and non-favourable outputs are discouraged or 'punished'.

SML has the task of learning a function that maps an input variable, like a set of features, to an output or label based on a set of input-output pairs used as examples. It is the most popular for performing ML operations. The algorithm works under supervision since it is provided with the actual outcome for each of the training inputs. The dataset is labelled, meaning that the data come with features and labels and the algorithm will carry out predictions or classification accordingly. As the training process goes on, the algorithm improves its capacity to identify the relationships between the two variables such that it could predict a new outcome. By supplying the supervised algorithm with more and more instances, it becomes able to learn more properly and predict an output more accurately. The success of the classification can be measured by testing the created model with a separate set of examples for which the true classifications are known but are hidden to the classifier. Intuitively, to evaluate the performance of different methods, a common practice is to divide the set of input-outputs pairs into two sets: training set, and test set. The training set is used to build the classifier model, while the test set is used to measure the accuracy of the classifier, i.e., it is a measure on how well it generalizes to unseen instances. Supervised learning could be distinguished in two main categories: predictive or directed so as well divided into two branches: Classification and Regression. Classification predictive modelling problems are different from regression predictive modelling problems. Classification is the task of predicting a

discrete class label. Regression is the task of predicting a continuous quantity. This part of the thesis uses SML algorithms for two main types of classification tasks:

1. Multi-Class Classification: it refers to those classification tasks that have more than two class labels. The examples are classified as belonging to one among a range of known classes.
2. Binary Classification: it refers to those classification tasks that have two possible class labels.

In this part of the thesis the features used as input for the SML algorithms are the data acquired by the iHealth Wave and they are described in the paragraph 3.3. Instead, the labels that are going to be predicted are the PH and Mind index that are transformed from numerical to categorical values as described in the following paragraphs. The most widely used learning algorithms for the classification include SVM, Linear and Logistic Regression, Naïve Bayes, DT, K-Nearest Neighbours algorithm (KNN), Neural Network (Multi-Layer Perceptron (MLP)), and RF. This is just a short list of the available algorithms. In the next paragraphs there will be a description of the ones adopted in this part of the thesis. The following Table 9 reports an example of the dataset used for the SML analysis.

FEATURES				LABELS	
Time in Bed (min)	Sleep Efficiency (%)	Nr. of Steps	Distance (km)	PH index	Mind index
500	95	1706.0	1.28	NORMAL	WELL
595	90	4380.0	3.31	NORMAL	WELL
540	94	3185.0	2.41	NORMAL	WELL
...	...	...	...	...	...
635	83	2622.0	1.97	BAD	BAD
...	...	...	...	...	...

*Table 9: example of dataset used for the supervised machine learning analysis.*

#### 4.3.3.1 Multi-class classification

To perform the Multi-Class Classification using the SML it is necessary to transform each numerical value of the PH, and Mind index into a categorical value whose possible labels are 'BAD', 'NORMAL' or 'WELL'. The conversion from numerical to categorical labels is accomplished according to the following criteria:

- if the PH or Mind index is lower than or equal to  $1/3$ , then the category 'Bad' is assigned to the PH or Mind index.

$$PH \text{ or Mind Index} \leq \frac{1}{3} \rightarrow PH \text{ or Mind Index} = 'Bad'$$

- if the PH or Mind index is greater than  $1/3$  and lower than or equal to  $2/3$ , then the category 'Normal' is assigned to the PH or Mind index.

$$PH \text{ or Mind Index} > \frac{1}{3} \text{ and } \leq \frac{2}{3} \rightarrow PH \text{ or Mind Index} = 'Normal'$$

- if the PH or Mind index is greater than  $2/3$ , then the category 'Well' is assigned to the PH or Mind index.

$$PH \text{ or Mind Index} > \frac{2}{3} \rightarrow PH \text{ or Mind Index} = 'Well'$$

The categorical PH and Mind index are used as labels for the classification algorithms.

The multi-class classification is further divided into 2 different approaches:

1. the multi-class classification without standardization of the features
2. the multi-class classification with standardization of the features. The main idea is to normalize (i.e.,  $\mu=0$  and  $\sigma=1$ ) each column of features individually, before applying the ML model.

#### 4.3.3.2 Binary classification

To perform the Binary Classification using the SML it is necessary to transform each numerical value of the PH, and Mind index into a categorical value whose possible labels are 'BAD', or 'NORMAL'. The conversion from numerical to categorical labels is accomplished according to the following criteria:

- if the PH or Mind index is lower than or equal to  $1/2$ , then the category 'Bad' is assigned to the PH or Mind index.

$$PH \text{ or Mind Index} < \frac{1}{2} \rightarrow PH \text{ or Mind Index} = 'Bad'$$

- if the PH or Mind index is greater than  $1/2$ , then the category 'Normal' is assigned to the PH or Mind index.

$$PH \text{ or Mind Index} \geq \frac{1}{2} \rightarrow PH \text{ or Mind Index} = 'Normal'$$

The categorical PH and Mind index are used as labels for the classification algorithms.

The multi-class classification is further divided into 2 different approaches:

1. the binary classification without standardization of the features
2. the binary classification with standardization of the features. The main idea is to normalize (i.e.,  $\mu=0$  and  $\sigma=1$ ) each column of features individually, before applying the ML model.

#### **4.3.3.3 Support Vector Machine**

In ML, the SVMs are supervised learning models with associated learning algorithms that analyse data for classification and regression analysis. Developed by Vladimir Vapnik with colleagues, SVMs are one of the most robust prediction methods, being based on statistical learning frameworks or VC theory, (56). Given a set of training examples, each belonging to one of two categories, a SVM algorithm creates a model that predicts if a new example falls to one class or the other, making it a non-probabilistic binary linear classifier (i.e., it predicts, for each given input, which of two possible classes the input is a member of). Intuitively, an SVM model represents the examples as points in space. The points are mapped so that the elements of the separate categories are divided by a clear gap that is as wide as possible. New examples are then mapped into that same space and predicted to fall into a category based on which side of the gap they belong to. More explicitly, a support-vector machine builds a hyperplane or set of hyperplanes in a high- or infinite-dimensional space, which can be used for classification, regression, or other tasks, (57). The hyperplane with the largest distance to the nearest training-data point of any class, so the maximum margin, allows to achieve a good separation, since in general the larger the margin, the lower the generalization error of the classifier, (58). The hyperplane dimensionality is equal to the number of input features minus one (i.e., with three feature the hyperplane will be a two-dimensional plane). These nearest training-data points are called Support Vectors (Figure 18). Support vectors are important because they are the training points that define the maximum margin of the hyperplane to the data set, and they therefore

determine the shape of the hyperplane. If one of them is moved and the SVM is retrained, the resulting hyperplane would change. So, SVM is an algorithm that takes the data as an input and outputs a line or hyperplane that separates the data into classes if possible. In addition to performing linear classification, SVMs can also carry out a non-linear classification using what is called the Kernel Trick, implicitly mapping their inputs into high-dimensional feature spaces (Figure 19), (59).

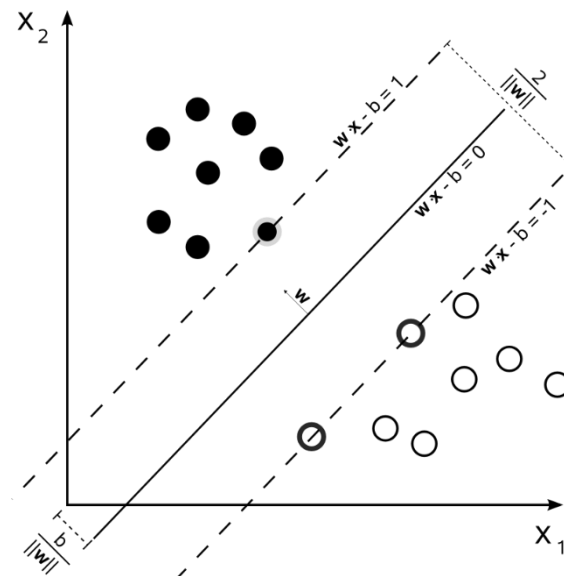


Figure 18: basic idea of SVM classifier.

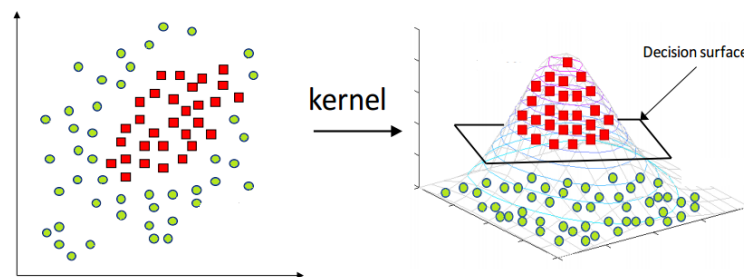


Figure 19: the training points mapped to a 3-D space where a separating hyperplane can be found.

If the input data are mapped into a higher-dimensional space, a linear algorithm operating in this space will behave non-linearly in the original input space. The Kernel trick is interesting because that mapping does not need to be ever computed. If the algorithm can be expressed only in terms of an inner product between two vectors, it is just necessary to replace this inner product with the inner product from some other suitable space. The kernel function accepts inputs in the original lower dimensional

space and returns the dot product of the transformed vectors in the higher dimensional space. That is the “trick”: every time a dot product is used, it is replaced with a Kernel function. Using the Kernel function, the algorithm can then be moved into a higher-dimension space without explicitly mapping the input points into this space. This is highly desirable, as sometimes the higher-dimensional feature space could even be infinite-dimensional and thus unfeasible to compute. Below is a list of the kernel functions used for the SML analysis:

1. the Linear Kernel is the simplest kernel function. Kernel algorithms using a linear kernel are often equivalent to their non-kernel counterparts. It is used when the data is Linearly separable, so using a single Line. It is mostly used when there are a large number of features in a particular Data Set (i.e., use Linear Kernel in Text Classification).
2. the Polynomial Kernel is a kernel function that represents the similarity of vectors (training samples) in a feature space over polynomials of the original variables, allowing learning of non-linear models. Intuitively, the polynomial kernel looks not only at the given features of input samples to determine their similarity, but also combinations of these. It is quite popular in natural language processing (NLP).
3. The Radial Basis Function (RBF) kernel is the most generalized form of kernelization and is one of the most widely used kernel due to its similarity to the Gaussian distribution.

#### **4.3.3.4 Decision Tree**

DT learning is one of the predictive modelling approaches used in statistics, data mining and ML. It uses a DT as a predictive model to move from features about an item, represented in the branches, to labels of the item, represented in the leaves, [(60), (61)]. The tree models where the labels can assume a discrete set of values are called classification trees: the leaves represent class labels and branches represent conjunctions of features that lead to those class labels. In decision analysis, a DT can be used to represent decisions and decision making graphically and explicitly. The goal is to create a model that predicts the value of a label by learning simple decision rules

inferred from the data features. Each internal (non-leaf) node of the DT is marked with an input feature. The arcs coming from a node marked with an input feature are labelled with each of the possible values of the target feature or the arc leads to a subordinate decision node on a different input feature. Each leaf of the tree is branded with a class, meaning that the data set has been classified by the tree into a specific class. The splitting is based on a set of splitting rules based on classification features. This process is repeated on each obtained subset in a recursive manner called recursive partitioning. The recursion is finished when the target variable of the subset at a node has all the same values, or when splitting no longer adds value to the predictions. The key terminologies related to DT are summarized below (Figure 20):

- Parent and Child Node: a node that gets divided into sub-nodes is known as Parent Node, while sub-nodes are called Child Nodes. A node can be divided into multiple sub-nodes; therefore, a node can act as a parent node of numerous child nodes.
- Root Node: the top-most node of a decision tree. It does not have any parent node. It represents the entire population or sample.
- Leaf/Terminal Nodes: nodes that do not have child nodes are known as Terminal/Leaf Nodes.

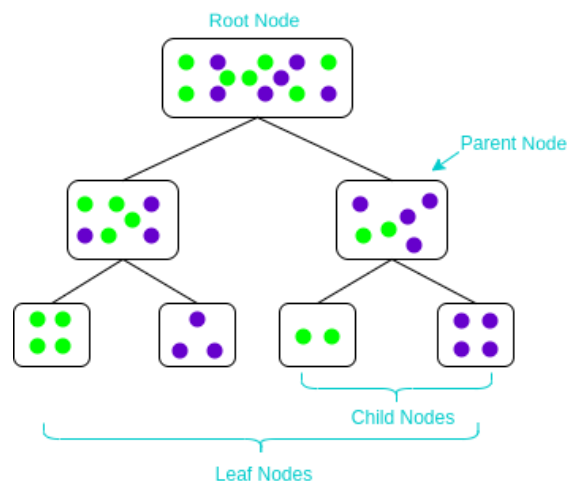


Figure 20: Decision Tree scheme.

Algorithms for build-up the DT usually work top-down, by choosing a variable at each step that best splits the set of features, (62). There are multiple ways of doing this, which can be broadly divided into two categories based on the type of target variable:

1. Continuous Target Variable
  - Reduction in Variance
2. Categorical Target Variable
  - Gini Impurity
  - Information Gain
  - Chi-Square

These metrics are applied to each candidate subset, and the resulting values are combined (e.g., averaged) to provide a measure of the quality of the split. Amongst other data mining methods, DTs have various advantages: they are simple to understand and interpret, they can handle both numerical and categorical data, they require little data preparation, they use a white/open box model, they have an in-built features selection, etc. however, they have also some limitations: DT can be very non-robust, and DT learners can create over-complex tree that do not generalize well from the training data.

#### **4.3.3.5 Random Forest**

RFs are an ensemble learning method for classification, regression and other tasks that work by building-up a multiple of DTs at training time. In the classification context, the output of the RF is the class selected by most of the trees (Figure 21). In the case of regression, the output is the mean or average prediction value of the individual DT, [ (63), (64), (65)]. RFs allow to correct the DTs' nature of overfitting to their training set. This ensemble of DTs is usually trained with the "bagging" method. The general idea of the bagging method is that a combination of learning models increases the overall result. The fundamental concept behind RF is the concept of wisdom of crowds. In data science speak, the RF model works well because a large number of relatively uncorrelated models (DTs) operating as an ensemble will outperform any of the individual founding models. RFs add additional randomness to the model, while building up the DTs. Instead of looking for the most important feature while splitting a node, it searches for the best feature among a random subset of features. This results in a wide diversity that generally results in a better model. Therefore, the RF considers only a random subset of the features for splitting a node. Another important characteristic of



the RF algorithm is that it is possible to measure the relative importance of each feature on the prediction very easily. Despite the RFs often achieve higher accuracy than a single DT, they however sacrifice the intrinsic interpretability of the DT. So, the RF are made of DTs that are not only trained on different sets of data (thanks to bagging) but also use different features to make decisions.

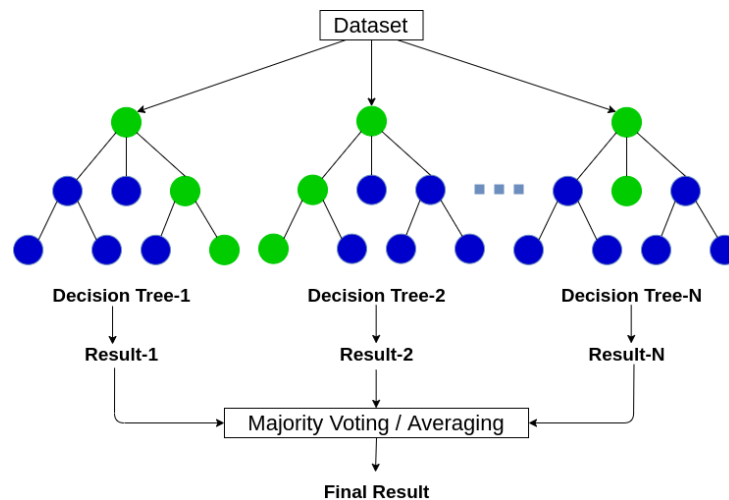
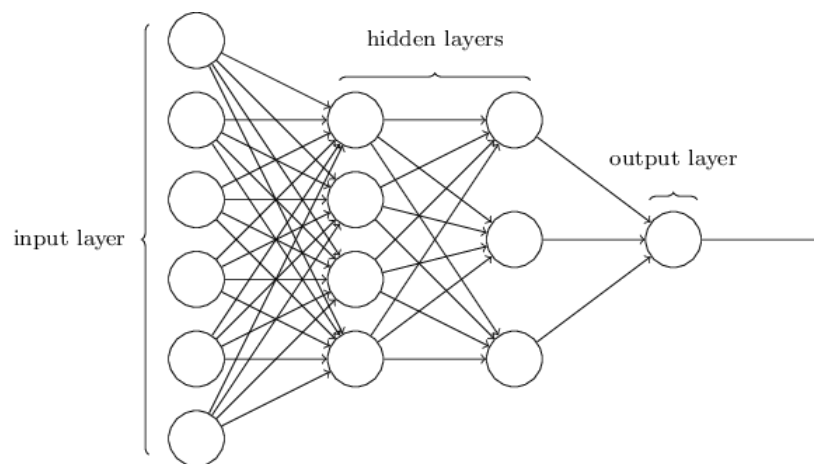


Figure 21: Random Forest Scheme

#### 4.3.3.6 MultiLayer Perceptron

A MultiLayer Perceptron (MLP) is a class of feedforward ANN. MLPs are sometimes colloquially referred to as "vanilla" neural networks, especially when they have a single hidden layer, (58). Perceptrons are inspired by the human brain and try to simulate its functionality to solve problems. An MLP consists of at least three layers of nodes: an input layer that receives the signal, a hidden layer that is the computational engine, and an output layer that makes a decision or prediction about the input (Figure 22). Except for the input nodes, each node is a neuron that uses a nonlinear activation function. MLP uses a supervised learning technique called backpropagation for training, (66). So, the correct label is compared with the output in order to compute the value of some predefined error-function. The error is fed back through the network and, by using this information, the MLP adjusts the weights of each connection. In this way it is possible to reduce the error function's value by some small amount. When the process is repeated for a sufficiently large number of training time, the network will usually converge to some state with a minimum error of the calculations. There are many non-

linear optimization methods used to adjust weights of the MLP properly. The optimizers define how to change the weights or learning rate to reduce the losses and provide the best results possible. The most used ones are: the Gradient Descent, the Stochastic Gradient Descent, the ADAM (67), etc. Its multiple layers and non-linear activation distinguish MLP from a linear perceptron. It can distinguish data that is not linearly separable. MLPs are trained on a set of input-output pairs and learn to model the correlation between those inputs and outputs. Training involves adjusting the parameters, or the weights and biases, of the model such to minimize error. Backpropagation is used to make those weigh and bias adjustments relative to the error, and the error itself can be measured in a variety of ways, including by root mean squared error (RMSE). Intuitively, MLPs are linked to one motion that goes back and one motion that goes forth. In the forward phase, the signal travels from the input layer towards the output layer passing through the hidden layers, while the decision of the output layer is established in regard to the ground truth labels. In the backwards phase, the set of weights and biases are backpropagated through the MLP.



*Figure 22: MultiLayer Perceptron scheme.*

#### **4.3.4 Performance metrics**

When the different models of SML are trained, getting the outputs in form of a label, the next phase is to evaluate their effectiveness using the performance metrics, (68). The evaluation metrics for binary classification models are Accuracy, Precision, Recall (or Sensitivity), and F1 Score. Before explaining in detail these metrics, it is important to introduce and explain the confusion matrix, also known as an error matrix, of the

evaluation of a binary or multi-class classification, (69). It is a tabular way to visualize the performance of a SML algorithm. Each row represents the instances in an actual class while each column represents the instances in a predicted class, or vice versa. Class labels in a binary classification problem can assume only two classes, preferably a positive and a negative class. The number of predictions where the classifier correctly predicts the positive class as positive and the negative class as negative are defined as True Positive (TP) value and True Negative (TN) value, respectively. Similarly, the number of predictions where the classifier incorrectly predicts the positive class as negative and the negative class as positive are defined False Negative (FN) value and False Positive (FP) value, respectively. The confusion matrix is simply a table that shows the number of elements within each of these four categories, as show in Figure 23 and 24. In contrast, in a typical multi-class classification problem, it is required to classify each sample into 1 of N different classes. So, unlike binary classification, there are no positive or negative classes. The definition for the TP is the same as in the binary confusion matrix. However, the TP is calculated for each class in the confusion matrix unlike the absolute TP of the binary case. The TN for a particular class is calculated by taking the sum of the elements in every row and column except the row and column of the class for which the TN is computed for. The TN for the other classes is computed alike. The FP for a particular class is computed by summing all the values in the column corresponding to that class but excluding the TP value. Consistently, the FN for a particular class is computed by summing all the values in the row corresponding to that class but excluding the TP value.

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Figure 23: binary Confusion Matrix.

		Predicted Class		
		Class 1	Class 2	Class 3
True Class	Class 1	TP	Class 2 predicted as Class 1	Class 3 predicted as Class 1
	Class 2	Class 1 predicted as Class 2	TP	Class 3 predicted as Class 2
	Class 3	Class 1 predicted as Class 3	Class 2 predicted as Class 3	TP

Figure 24: example of 3-class Confusion Matrix.

The classification Accuracy both for binary and multi-class classification is the ratio of correct prediction to the total number of input samples, so it is the percentage of instances correctly classified. The formula of the classification accuracy is:

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} \cdot 100 \quad (4)$$

The Precision is the fraction of relevant instances among the retrieved instances. Its values range from 1 (optimum) to zero. The formula of the Precision is:

$$Precision = \frac{TP}{TP + FP} \quad (5)$$

The Recall or Sensitivity quantifies the number of correct positive predictions made out of all positive predictions that could have been made. The best value is 1 and the worst value is 0. It provides an indication of missed positive predictions. The formula of the Recall is:

$$Recall = \frac{TP}{TP + FN} \quad (6)$$

The F1-Score can be interpreted as a weighted average of the Precision and Recall, where and the score reaches its best value at 1 and worst at 0. The relative contribution of Precision and Recall to the F1-Score are equal. The formula of the F1-Score is:

$$F1 - Score = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall} \quad (7)$$

The Leave-One-Out Cross-Validation (LOOCV) is used to validate the SML algorithms implemented in the thesis. LOOCV is an extreme configuration of k-fold cross-validation where  $k$  is set to the number of examples in the dataset. The K-fold cross-validation process consists in the partitioning the dataset in  $k$  parts of equal size. One part is used as a test set, while the remaining  $k-1$  parts are used as training sets. The process is iterated  $k$  times, where at each iteration each part of the dataset is used exactly once as a test set.

## 4.4 Results

The analysis described in the previous chapter is applied to the data of the User 1. User 2 has been excluded since his questionnaire data are homogeneous and therefore not useful for the training of the algorithms. Figure 25 shows accuracies of the SML algorithms used for the 3-class classification of the PH index. Figure 26 reports the same score for the Mind index. The other performance metrics are represented together with different colours in the Figure 27 and Figure 28 for the PH and Mind index, respectively. The accuracy and the performance metrics of the remaining SML analysis are summarized in the Table 10 and in the Table 11 and it is possible to see that, higher accuracy for the prediction of the PH index is achieved by the SVM with Polynomial and RBF kernel in all the 4 types of analysis performed, with an accuracy of 98 and 100%. The goodness of the 2 models is also validated by the additional metrics provided in Table 12, 13, and 14, which are the Precision, the Recall, and the F1-score. These metrics are important because they provide a more complete picture of the effective classification results regarding the goodness of the model in predicting the output quantity; the range of this metrics for the trained algorithms is between 0.64 and 1.0, so the predictions of the models are mostly correct, making the majority of the input features corresponding to the right class or output labels.

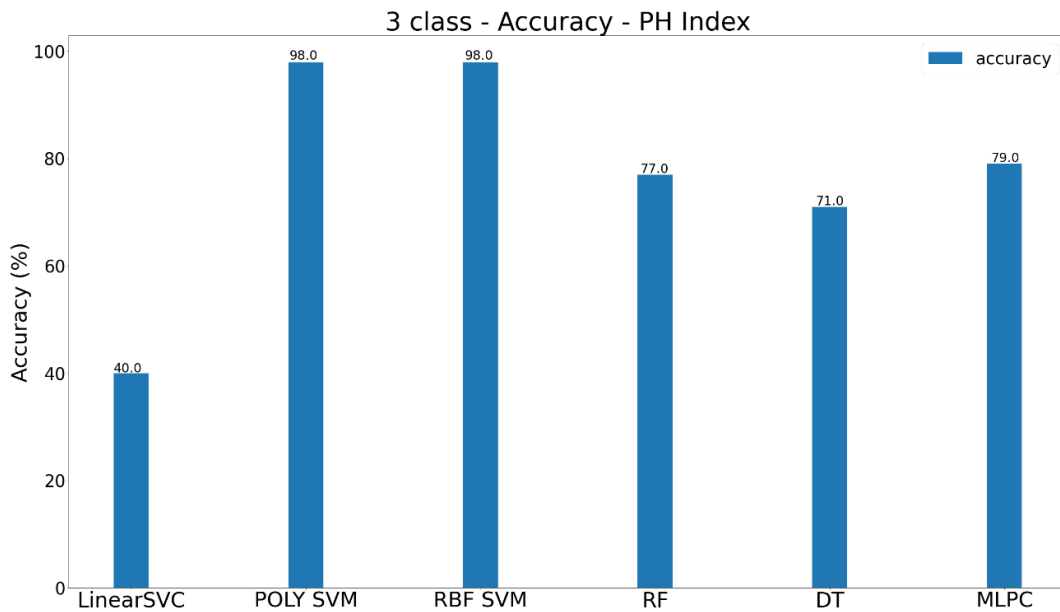


Figure 25: Accuracy score for the 3-class prediction of the PH index for the User 1.

Typology	ACCURACY (%) – PH index					
	Linear SVC	Poly SVM	RBF SVM	RF	DT	MLPC
3 class	40 %	98 %	98 %	77 %	71 %	79 %
3 class + StS	79 %	98 %	98 %	77 %	70 %	70 %
2 class	59 %	100 %	100 %	80 %	68 %	72 %
2 class + StS	52 %	100 %	100 %	81 %	67 %	73 %

Table 10: Accuracy score for the prediction of the PH index for the User 1.

Regarding the Mind index, the SVM with RBF kernel achieves the best accuracy of 100% and the best score of 1 for the 3 additional performance metrics in the case of the 3-class classification. On the contrary, in the case of the 3-classification with the Standard Scaler (StS) function, the performances of the SVM with RBF kernel slightly decrease, achieving an Accuracy of 78% while the Precision, Recall, and F1-Score vary in the range between 0.55 and 0.62. The performance of the model increases again, achieving the best possible values, in both cases of the binary classification.

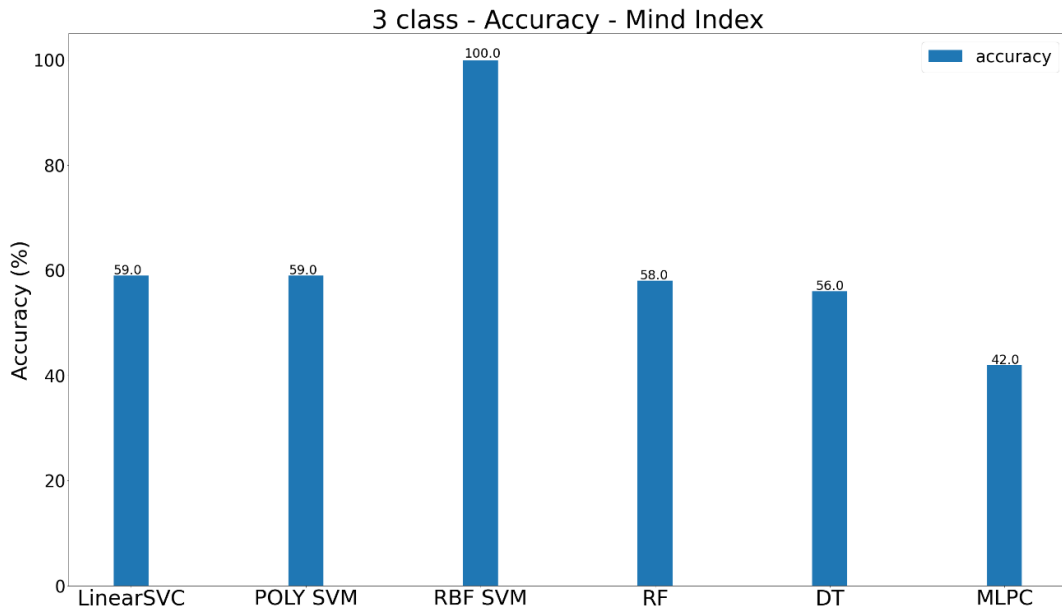


Figure 26: Accuracy score for the 3-class prediction of the Mind index for the User 1.

Typology	ACCURACY (%) – Mind index					
	Linear SVC	Poly SVM	RBF SVM	RF	DT	MLPC
3 class	59 %	59 %	100 %	58 %	56 %	42 %
3 class + Sts	73 %	62 %	78 %	59 %	55 %	63 %
2 class	44 %	85 %	100 %	83 %	77 %	78 %
2 class + Sts	55 %	100 %	100 %	83 %	76 %	85 %

Table 11: Accuracy score for the prediction of the Mind index for the User 1.

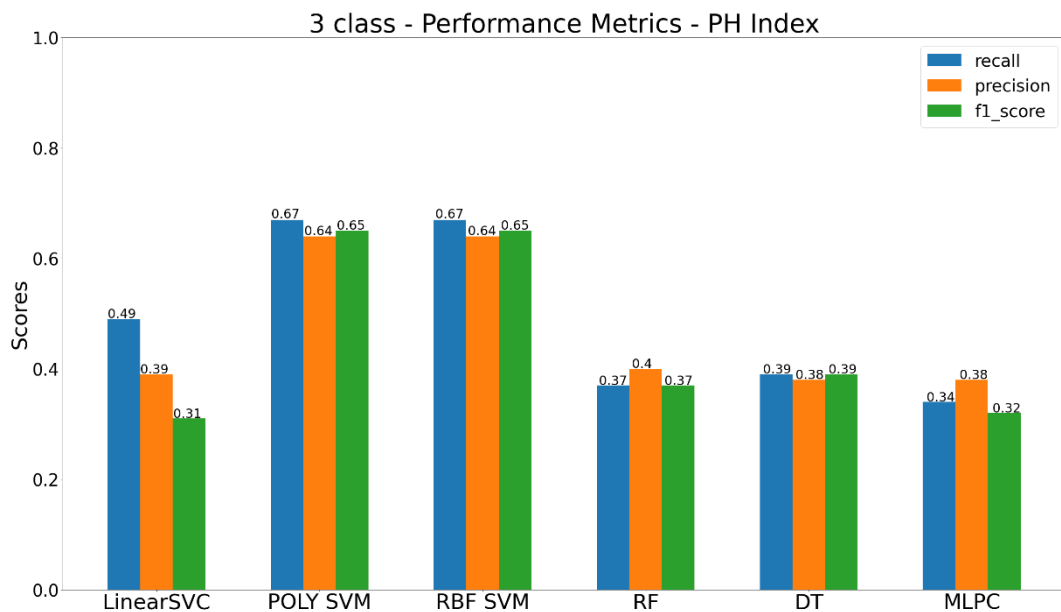


Figure 27: performance metrics for the 3-class prediction of the PH index for User 1.



Figure 28: performance metrics for the 3-class prediction of the Mind index for User 1.

Typology	PRECISION											
	Linear SVC		Poly SVM		RBF SVM		RF		DT		MLPC	
	PH	MIND	PH	MIND	PH	MIND	PH	MIND	PH	MIND	PH	MIND
3 class	0.39	0.2	0.64	0.2	0.64	1.0	0.4	0.47	0.38	0.49	0.38	0.32
3 class + StS	0.26	0.52	0.64	0.52	0.66	0.55	0.4	0.47	0.39	0.49	0.46	0.57
2 class	0.44	0.59	1.0	0.43	1.0	1.0	0.7	0.53	0.53	0.58	0.52	0.55
2 class + StS	0.57	0.53	1.0	1.0	1.0	1.0	0.72	0.53	0.56	0.57	0.57	0.69

Table 12: Precision metrics for the prediction of the PH and Mind index for User 1.

Typology	RECALL											
	Linear SVC		Poly SVM		RBF SVM		RF		DT		MLPC	
	PH	MIND	PH	MIND	PH	MIND	PH	MIND	PH	MIND	PH	MIND
3 class	0.49	0.33	0.67	0.33	0.67	1.0	0.37	0.42	0.39	0.49	0.34	0.32
3 class +	0.33	0.53	0.67	0.46	0.67	0.62	0.37	0.42	0.4	0.49	0.52	0.47
2 class	0.43	0.65	1.0	0.5	1.0	1.0	0.62	0.51	0.54	0.59	0.53	0.55
2 class +	0.6	0.55	1.0	1.0	1.0	1.0	0.62	0.51	0.56	0.58	0.55	0.59

Table 13: Recall metrics for the prediction of the PH and Mind index for User 1.

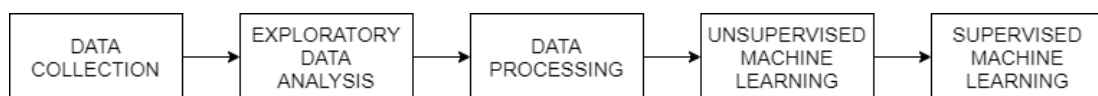
Typology	F1 Score											
	Linear SVC		Poly SVM		RBF SVM		RF		DT		MLPC	
	PH	MIND	PH	MIND	PH	MIND	PH	MIND	PH	MIND	PH	MIND
3 class	0.31	0.25	0.65	0.25	0.65	1.0	0.37	0.43	0.39	0.4	0.32	0.32
3 class + StS	0.29	0.52	0.65	0.47	0.66	0.57	0.37	0.42	0.39	0.49	0.27	0.3
2 class	0.43	0.42	1.0	0.46	1.0	1.0	0.63	0.5	0.54	0.58	0.52	0.55
2 class + StS	0.5	0.47	1.0	1.0	1.0	1.0	0.64	0.5	0.57	0.58	0.55	0.61

Table 14: F1 Score metrics for the prediction of the PH and Mind index for User 1.



# Chapter 5: Unsupervised Machine Learning analysis of elderlies' behavioral data

The main objective of the third part of the thesis is to propose a methodology that can analyse the domotic data acquired using a network of PIR and Light sensors installed in a set of multi or single resident private homes in order to profile the users' daily behaviours and habits. The profile of the users can be identified using the Unsupervised ML (UML) algorithm called K-Means. The algorithm has been implemented in Python and the work took place in the workplace of Spyder. The Methodology adopted in the third part of the thesis can be divided in 5 steps as follows:



## 5.1 Domotic data collection

The Domotic Data used in the third part of the thesis are acquired using a wireless sensor network composed by a variable number of PIR and Light sensors that were installed in the most important rooms and furniture of the house. Each private house had a different sensing framework. Figure 29 displays an example of one house's planimetry with its relative sensors installed in the most informative rooms. The time window of the data was variable among the different dwellings, but in general it was included between February 2018 and May 2019. The data collected consists of a daily list of PIR and Light sensors' activations with their relative Sensor Status and Timestamp. The PIR's sensor status can assume only the value '1' when it is activated, while the sensor status of the Light can assume a binary value, '1' or '0', depending on if it is switched On or Off. The timestamp gives the Date and Time of day associated with PIR activation or Light activation and deactivation. The participants involved in the study live in 7 different homes that are single-resident and multi-resident apartments. Their age is between 68-82 and they did not have particular pathologies. This group of participants was living alone or in couple, with normal social status, no psychological disease and able to perform everyday activities. Table 15 reports an example of domotic dataset for the House nr. 3.



Figure 29: schematic plan of one of the apartments with installed domotic sensors.

User_Id	User	Device_Id	Time_Stamp	Feature_Id	Feature_Name	Value
5a0563...	home_3	5a0567...	01/02/2018 19:50:43	5a0b0...	PIR_room1	1
5a0563...	home_3	5a0567...	01/02/2018 19:52:58	5a0b0...	PIR_room1	1
5a0563...	home_3	5a0567...	01/02/2018 19:53:07	5a0b0...	Light_6	1
5a0563...	home_3	5a0567...	01/02/2018 19:53:12	5a0b0...	PIR_room2	1
5a0563...	home_3	5a0567...	01/02/2018 19:53:31	5a0b0...	Light_6	0
...	...	...	...	...	...	...

Table 15: example of the Domotic Data of the Home nr. 3.

## 5.2 Exploratory data analysis

After the Data Collection, the following step is the removal of all the possible outliers. Primarily, the first and last day of each house's dataset are removed because they are not complete. After that, it is necessary to delete from the data all the days that do not contain any sensors activations or the ones with at least one sensor's malfunctioning. The light sensors' malfunctions may include the case of absence of one of the 2 possible sensor status following the sensor's opposite condition. Intuitively, the Light status '1' or '0' is not followed by the opposite state, but by the same condition.

## 5.3 Data processing

From the general Dataset it is possible to extract 2 subsets containing one the PIR sensors and the other only the Light sensors. The following step is the splitting of each day of the 2 subsets in 6 time slots of 4 hours each:

TIME SLOT 1	TIME SLOT 2	TIME SLOT 3	TIME SLOT 4	TIME SLOT 5	TIME SLOT 6
12 - 04 am	04 – 08 am	08 am – 12 pm	12 – 04 pm	04 – 08 pm	08 pm – 12 am

Within each time slots it is necessary to compute the sum of the activations of each sensor. Subsequently, the total number of activations for both types of equipment are saved into 6 different datasets, one for each of the time slots. Intuitively, the table contains a number of columns equal to the sum of the PIR and Light sensors deployed in each home and a number of rows equal to the number of days of the general Dataset. Each element of one row of a dataset contains the sum of activations of the corresponding sensor (column) inside the same time slot and for a single day of the dataset. The following Table 16 reports an example of data for the first time slot.

Date	Time_slot	light1	light3	light4	light5	light6	light8	light9	PIR1	PIR2
03/02/2018	12 – 04 am	0	0	1	5	1	2	2	6	5
04/02/2018	12 – 04 am	0	0	0	1	0	0	0	0	2
05/02/2018	12 – 04 am	2	0	2	4	1	2	2	8	3
06/02/2018	12 – 04 am	1	0	1	3	0	1	1	2	7
...	...	...	...	...	...	...	...	...	...	...

Table 16: example of Table for Time Slot 12 – 04 am of the Home nr. 3.

## 5.4 Unsupervised Machine Learning

UML algorithms infer patterns from a dataset without knowing neither the target labels nor the internal structure. Unlike SML, UML methods cannot be directly employed to a regression and a classification problem because the values for the output data are unknown, making it impossible the training of the algorithms. UML can instead be used to discover any similarities, differences, patterns, and underlying structure in the data. A dataset is provided to the UML model. The algorithm analyses the data and eventually finds a classification criterion. It is the model itself that finds the labels to be assigned to the examples, (70). In UML, the 2 most used methods are:

1. Clustering: it is used in UML to group, or divide, datasets with common attributes in order to infer algorithmic relationships.
2. Dimensionality Reduction: it is used to reduce the number of data inputs (features or dimension) to a manageable size while also preserve the integrity of the dataset.

### 5.4.1 Clustering Analysis

Cluster Analysis or Clustering is a data mining technique grouping unlabelled data based on their similarities or differences. The clustering algorithms process raw and unclassified data elements into groups, called cluster, represented by structures or patterns in the information. The data points are divided into a number of groups such that the elements in the same groups are more similar to other ones in the same group than those in other groups. There are various algorithms to solve the clustering task that differ basically in their way of thinking of what constitutes a cluster and how to efficiently discover them. Even if the notion of a "cluster" is subjective and not exactly defined, the most popular include groups with small distances between cluster members, dense areas of the data space, intervals, or particular statistical distributions. Therefore, there are many methods for achieving the same goal. Every methodology adopts a specific set of rules for defining the 'similarity' among data elements. There are more than 100 clustering algorithms, but few of the algorithms are more popular, [(71), (72)]:

- Connectivity models are based on the concept that the data points closer in data space are more similar to each other than the data points lying farther away. Hierarchical clustering algorithm and its variants are example of these models.
- Centroid models that are iterative clustering algorithms where the similarity is inferred by the proximity of an element to the centroid of the clusters. K-Means clustering algorithm is a popular algorithm that falls into this category.
- Distribution models compute the probability that a point belongs to a cluster. The algorithm defines, around each possible centroid, the density distributions for each cluster, measuring the probability of belonging based on those

distributions. The Expectation-maximization algorithm is an example of these models.

- Density Models search for areas of diverse density of points in the data space. It separates many density regions and assign the points within these area to the same cluster. DBSCAN and OPTICS are 2 popular examples.

All of these models can be further classified into four distinct categories, (73):

- Exclusive Clustering: the data are grouped in exclusive way, making each element belong to a single cluster.
- Overlapping Clustering: it uses fuzzy sets to cluster data, so each element can belong to more than one cluster with different membership's degree.
- Hierarchical Clustering: hierarchical structures are used to organize the clusters.
- Probability clustering, where a probabilistic approach is adopted.

#### **5.4.1.1 K-Means**

K-means clustering is one of the simplest and popular UML algorithms for clustering operations, (74). It is an Exclusive Clustering algorithm. Its aim is to divide the data into a fixed number of 'k' clusters such that the elements in the same cluster are similar and elements in the different clusters are farther apart. Similarity of two points is determined by the distance between them. Figure 30 summarize the steps of K-Means algorithm. Training examples are shown as dots, and cluster centroids are shown as crosses (Figure 30 (a)). The first step of the K-means is to define 'k' centroids, one for each cluster. These centroids should be placed in a smart way since different position cause different results. So, the clever choice is to place them as far away as possible from each other (Figure 30 (b)). The next step is to pick up each point belonging to the data set and associate it to the closest centroid. When no point is left, the first step is finished, and an initial grouping is done (Figure 30 (c)). At this point it is necessary to re-calculate 'k' new centroids as barycentre of the previous step's resulting clusters. When the 'k' new centroids are found, a new association has to be done between the same elements of the data set and the closest new centroid. These 2 steps are repeated in loop, with the 'k' centroids change their location step by step, until no more changes are done (Figure 30 (d-f)). In simple words, the 'k' centroids do not move any more, [ (75),

(76)]. Finally, the algorithm aims to divide these data point into 'k' cluster such to minimize the objective function called Within-Cluster Sum of Squares (WCSS). Even though it can be proved that the K-means algorithm will always end, it does not necessarily find the most optimal configuration. Moreover, the algorithm is also sensitive to the initial randomly selected centroids.

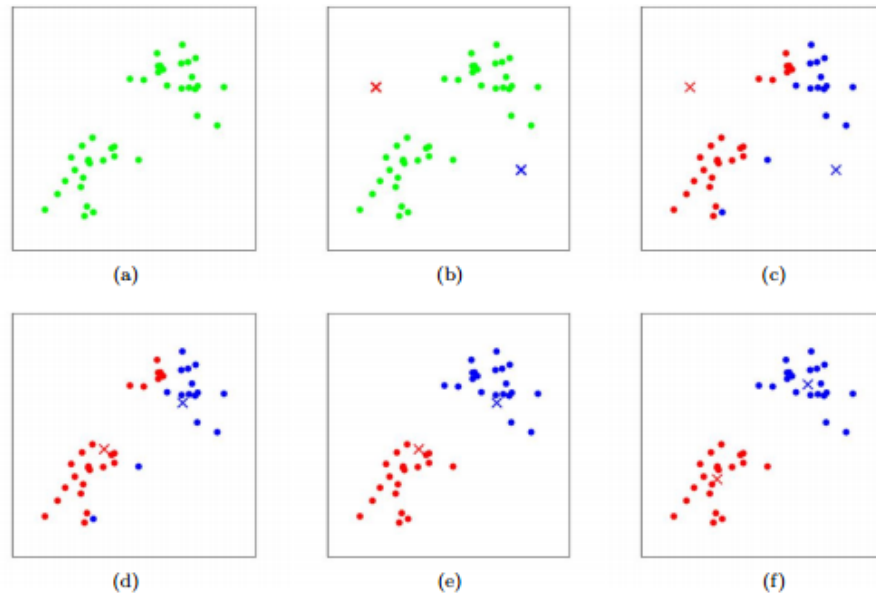


Figure 30: K-Means algorithm.

Each of the 6 datasets containing the daily number of activations for each sensor within the same time slot are used as input of the K-Means algorithm implemented in Python. The Principal Component Analysis (PCA) is applied to the 6 datasets before of the clustering algorithm. The PCA is a dimensionality-reduction method used to reduce the dimensionality of large data sets, by transforming a large set of variables into a smaller one, while preserving as much information as possible of the large set. The 'k' parameter of the algorithm is set to 2, meaning that each element of the tables can be associated with only 2 possible clusters.

#### 5.4.1.2 Principal Components Analysis

The PCA is a statistical procedure that performs an orthogonal linear transformation to convert a set of data of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components, (77). The number of principal components is usually less than the original number of variables. It is commonly used

for dimensionality reduction: each data point is projected onto only the first few principal components to get lower-dimensional data while preserving the data's variation (i.e., statistical information) as much as possible. Intuitively, the PCA transforms the data to a new coordinate system in such a way that the first principal component, with the largest possible variance, will lie on the first coordinate, the second greatest variance will lie on the second coordinate orthogonal to the preceding one, and so on. The 1<sup>st</sup> principal component can be defined as a direction that maximizes the variance of the projected data. The  $i$ -th principal component can be defined as a direction orthogonal to the previous  $i-1$  principal components that maximizes the variance of the projected data. Mathematically, the entire process of get the principal components from a high dimensional dataset can be summarized in six steps, (78) (Figure 31):

- 1) Take the whole dataset made of  $d+1$  dimensions and neglect the labels such that the new dataset is now  $d$  dimensional.
- 2) Compute the mean for every dimension of the whole dataset.
- 3) Compute the covariance matrix of the whole dataset.
- 4) Compute eigenvectors and the corresponding eigenvalues.
- 5) Sort the eigenvectors by decreasing eigenvalues and choose  $k$  eigenvectors with the largest eigenvalues to form a  $d \times k$  dimensional matrix  $W$ .
- 6) Use the  $d \times k$  eigenvector matrix to transform the samples onto the new subspace.

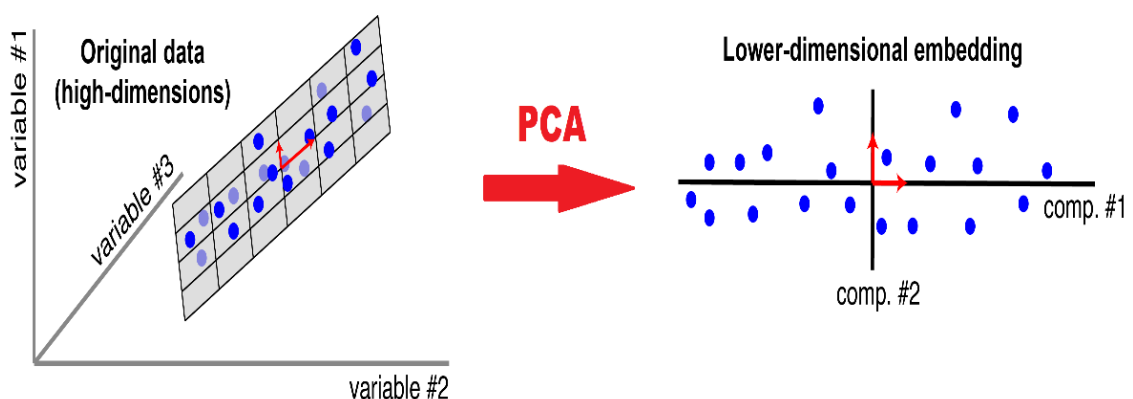


Figure 31: transformation of high dimensional data to low dimensional data via PCA.

## 5.4.2 Silhouette Score

The Silhouette score is used to evaluate the quality of clusters created by the K-Means algorithm in terms of how well data points are clustered with other data points that are similar to each other, (79). In order to calculate the Silhouette score for each observation/data point, it is necessary to compute the following distances for each observation belonging to all the clusters:

- Mean distance between the element and all other data points in the same cluster. This distance can also be called a mean intra-cluster distance and it is denoted by  $a$ .
- Mean distance between the element and all other data points of the next nearest cluster. This distance can also be called a mean nearest-cluster distance and it is denoted by  $b$ .

For each data point  $i$ , the Silhouette score is calculated using the following formula:

$$S_i = \frac{b_i - a_i}{\max [a_i, b_i]} \quad (8)$$

The value of the Silhouette score varies between -1 and 1. If the score is 1, the cluster is dense and well-separated from the other clusters. A value near 0 represents overlapping clusters where the samples are very close to the decision boundary of the neighbouring clusters. A negative score between -1 and 0 indicates that the samples might have got assigned to the wrong clusters. The Silhouette score for the 6 datasets used in the thesis is given as the mean of the Silhouette score for each of the 6 set of observations.

## 5.5 Clustering interpretation using Supervised Machine Learning

The last step of the analysis is the interpretation of the clustering results. This is done by using a set of rules automatically generated by training a supervised DT model employing the 6 tables with the number of activations as the features and the clustering results as labels, [ (80), (81)]. The DT model, after the training, will report the clustering



rules for the data separation in the different groups and the proportion of the data within the cluster satisfying that rule.

## 5.6 Results

The analysis described in the previous chapter is performed on each house’s dataset composed by 6 tables containing each the sensors’ number of activations within the same time slot. Table 17 reports the Silhouette score computed for the different time slots of each home. The coefficient is computed as the mean of the Silhouette scores for each element of the dataset considered. Highest silhouette scores are obtained for House 4 and House 5 for the first time slot, suggesting that the two clusters are dense and well separated, while the low Silhouette score of the House 7 in the 2<sup>nd</sup> time slot indicates that the two clusters are not clearly separable.

Time Slot	SILHOUETTE SCORE						
	House 3	House 4	House 5	House 6	House 7	House 8	House 9
12 – 04 am	0.486	0.659	0.852	0.563	0.525	0.356	0.496
04 – 08 am	0.426	0.472	0.432	0.404	0.342	0.545	0.632
08 am – 12 pm	0.409	0.560	0.418	0.446	0.366	0.442	0.623
12 – 04 pm	0.411	0.517	0.479	0.515	0.396	0.654	0.531
04 – 08 pm	0.437	0.476	0.461	0.520	0.411	0.459	0.601
08 pm – 12 am	0.445	0.527	0.404	0.443	0.524	0.402	0.545

*Table 17: Silhouette Score for the dataset.*

The clustering results are difficult to interpret without an explanation of separation rules that provide them. In this regard, the following tables report the clustering rules identified by the DT model for each time slot of the different houses. These tables help to understand which sensors are responsible for the assignment of ‘diversity’ condition of some days respect to the remaining ones. Therefore, each time slot of the different houses has a different clustering rules that is provided by a different PIR or Light sensor. As an example, Table 18 reports the results of the House 4 obtained from K-means clustering algorithms that are interpreted by the DT model. The clustering rules based on the most informative sensors are reported between round brackets for each time slot. The ratio of days included in one of the 2 clusters and respecting its clustering rules is reported between square brackets. The number outside of the brackets indicated the

number of days belonging to each cluster. The same explanation can be applied to the remaining tables reporting the results for the other houses.

Time slot	Cluster 1	Cluster 2
12 – 04 am	199 [0.99] (PIR2-Bedroom <= 1.5)	101 [1.0] (PIR2-Bedroom > 1.5)
04 – 08 am	162 [0.89] (PIR2-Bedroom <= 4.5)	138 [0.98] (PIR2-Bedroom > 4.5)
08 am – 12 pm	178 [0.97] (PIR2-Bedroom > 9.5)	122 [0.98] (PIR2-Bedroom <= 9.5)
12 – 04 pm	183 [0.93] (PIR1-Living Room <= 10.5)	117 [0.93] (PIR1-Living Room > 10.5)
04 – 08 pm	125 [0.89] (PIR1-Living Room > 9.5)	175 [0.93] (PIR1-Living Room <= 9.5)
08 pm – 12 am	185 [0.96] (PIR2-Bedroom > 3.5)	115 [0.99] (PIR2-Bedroom <= 3.5)

Table 18: clustering rules for the House 4.

Time slot	Cluster 1	Cluster 2
12 – 04 am	196 [0.88] (light5-Bathroom <= 2.5)	83 [0.78] (light5-Bathroom > 2.5)
04 – 08 am	159 [0.88] (PIR2-Bedroom <= 8.5)	120 [0.91] (PIR2-Bedroom > 8.5)
08 am – 12 pm	116 [0.89] (PIR1-Living Room <= 17.5)	163 [0.95] (PIR1-Living Room > 17.5)
12 – 04 pm	163 [0.99] (PIR1-Living Room > 15.5)	116 [0.97] (PIR1-Living Room <= 15.5)
04 – 08 pm	116 [0.94] (PIR1-Living Room <= 19.5)	163 [0.96] (PIR1-Living Room > 19.5)
08 pm – 12 am	109 [0.99] (PIR1-Living Room <= 14.5)	170 [0.99] (PIR1-Living Room > 14.5)

Table 19: clustering rules for the House 3.

Time slot	Cluster 1	Cluster 2
12 – 04 am	369 [1.0] (PIR1-Living Room <= 7.0)	12 [1.0] (PIR1-Living Room > 7.0)
04 – 08 am	208 [0.95] (PIR1-Living Room <= 6.5)	173 [0.79] (PIR1-Living Room > 6.5)
08 am – 12 pm	202 [0.98] (PIR1-Living Room <= 16.5)	179 [0.80] (PIR1-Living Room > 16.5)
12 – 04 pm	284 [0.96] (PIR1-Living Room <= 14.5)	97 [0.89] (PIR1-Living Room > 14.5)
04 – 08 pm	135 [0.95] (PIR1-Living Room > 16.5)	246 [0.95] (PIR1-Living Room <= 16.5)
08 pm – 12 am	244 [0.93] (PIR1-Living Room <= 10.5)	137 [0.91] (PIR1-Living Room > 10.5)

Table 20: clustering rules for the House 5.

Time slot	Cluster 1	Cluster 2
12 – 04 am	95 [1.0] (PIR2-Bedroom <= 1.5)	46 [1.0] (PIR2-Bedroom > 1.5)
04 – 08 am	67 [0.90] (PIR1-Living room <= 15.5)	74 [0.98] (PIR1-Living room > 15.5)
08 am – 12 pm	53 [0.97] (PIR1-Living room > 22.5)	88 [0.89] (PIR1-Living room <= 22.5)
12 – 04 pm	41 [1.0] (PIR1-Living room > 22.5)	100 [1.0] (PIR1-Living room <= 22.5)
04 – 08 pm	44 [1.0] (PIR1-Living room > 21.5)	97 [0.98] (PIR1-Living room <= 21.5)
08 pm – 12 am	85 [0.96] (PIR1-Living room <= 12.5)	56 [1.0] (PIR1-Living room > 12.5)

Table 21: clustering rules for the House 6.

Time slot	Cluster 1	Cluster 2
12 – 04 am	202 [0.99] (PIR3-Bathroom <= 2.5)	99 [1.0] (PIR3-Bathroom > 2.5)
04 – 08 am	135 [0.91] (PIR3-Bathroom > 10.5)	166 [0.93] (PIR3-Bathroom <= 10.5)
08 am – 12 pm	178 [0.85] (PIR3-Bathroom <= 22.5)	123 [0.94] (PIR3-Bathroom > 22.5)
12 – 04 pm	175 [0.89] (PIR3-Bathroom <= 13.5)	126 [0.92] (PIR3-Bathroom > 13.5)
04 – 08 pm	176 [0.85] (PIR3-Bathroom <= 10.5)	125 [0.85] (PIR3-Bathroom > 10.5)
08 pm – 12 am	69 [0.80] (PIR3-Bathroom > 5.5)	232 [0.96] (PIR3-Bathroom <= 5.5)

Table 22: clustering rules for the House 7.

Time slot	Cluster 1	Cluster 2
12 – 04 am	147 [0.91] (PIR1-Living room <= 3.5)	129 [0.84] (PIR1-Living room > 3.5)
04 – 08 am	158 [0.98] (PIR1-Living room > 11.5)	118 [0.99] (PIR1-Living room <= 11.5)
08 am – 12 pm	117 [0.94] (PIR1-Living room <= 23.5)	159 [0.98] (PIR1-Living room > 23.5)
12 – 04 pm	237 [1.0] (PIR1-Living room > 19.5)	39 [1.0] (PIR1-Living room <= 19.5)
04 – 08 pm	195 [0.98] (PIR1-Living room > 25.5)	81 [1.0] (PIR1-Living room <= 25.5)
08 pm – 12 am	102 [0.83] (PIR1-Living room <= 14.5)	174 [0.95] (PIR1-Living room > 14.5)

Table 23: clustering rules for the House 8.

Time slot	Cluster 1	Cluster 2
12 – 04 am	138 [0.95] (PIR2-Bedroom <= 3.5)	84 [0.92] (PIR2-Bedroom > 3.5)
04 – 08 am	182 [0.95] (PIR3-Bathroom <= 9.5)	40 [0.91] (PIR3-Bathroom > 9.5)
08 am – 12 pm	167 [0.97] (PIR1-Living room <= 39.5)	56 [1.0] (PIR1-Living room > 39.5)
12 – 04 pm	49 [0.96] (PIR1-Living room > 39.5)	173 [0.99] (PIR1-Living room <= 39.5)
04 – 08 pm	173 [0.98] (PIR1-Living room <= 35.5)	49 [0.94] (PIR1-Living room > 35.5)
08 pm – 12 am	171 [0.93] (PIR1-Living room <= 18.5)	51 [0.85] (PIR1-Living room > 18.5)

Table 24: clustering rules for the House 10.

As it has been mentioned in previous lines, each time slot of each house has a different clustering rules based on the most informative sensors that is used to separate the days between the 2 clusters.

## Chapter 6: Conclusions

In this thesis, the analysis presented has the final aim of improving elderlies' well-being using the data collected by non-intrusive sensors, such as wearable sensors and domotic sensors. Data are processed through traditional methods, like statistical analysis, and ML approaches, both supervised and unsupervised. Data analysed in the first two sections of this work are collected from a commercially available smartwatch worn by

older adults, while data collected in the third part are obtained through a domestic wireless sensors network, composed of PIR and Light sensors, installed in 7 different single or multi-resident homes inhabited by elderly, to track their behaviour at home. The results of the statistical analysis (chapter 3.5), show that by computing a 2-weeks sliding period over the smartwatch data (e.g., number of steps, sleep efficiency), is a reasonable time window to classify the following available day as 'Normal' or 'Abnormal'. The latter class can be further categorized as 'Negative' or 'Positive' abnormality depending on the comparison between the data available and the range of values considered. The results of this classification can be used as decision criteria to provide different coaching solutions to the elderly user, to maintain an active and healthy lifestyle and to improve sleep habits. The coaching system can congratulate in case of positive abnormality, suggest improvements when the previous day is in the normality range, or asks the elderly if a medical examination is necessary to solve a physical or mental issue in the last case. It is also necessary to consider that the smartwatch's measurement uncertainty on the daily steps count of 4.7% may be responsible of a wrong coaching as demonstrated by 6.7% and 11.4% uncertain days out of the total for the 2 users considered in the analysis. Chapter 4.4, which shows results related to SML for predicting elderly's well-being, concludes that the physical and mental wellness, summarized by 2 indices computed from daily questionnaires' answers, can be predicted by training the SML algorithms using the smartwatch data as features and as labels the categorical values summing up the wellness condition of the elderly reported in the surveys. The various types of SML analysis performed have demonstrated the superior performances of the SVM with RBF kernel for predicting the PH and Mind indices. Indeed, the accuracy of the algorithm is included between 78% and 100%. The last part of the thesis demonstrates that the combination of UML and SML can be used to profile the elderly living in a smart home equipped with PIR and Light sensors by separating the daily number of activations of sensors in 6 time slots among the different domestic environments. K-means clustering algorithms has been applied on each dataset, to detect any possible deviation in the sensors' activation during the 6 time slots of daily living of the elderly. Best results are obtained for the 1<sup>st</sup> time slot of the House 4 and House 5 with a Silhouette score of 0.659 and 0.852,

suggesting clear separation between the days belonging to 2 different clusters. Moreover, DT model used after the UML analysis of each time slot helps to understand the clustering methodology adopted by the K-means algorithm for each time window of the different houses. This makes possible to identify which sensor may cause the assignment of a 'diversity' condition between the days of the dataset. Each time slot is characterized by a specific clustering rule that is provided by a different sensor. Based on the defined separation rules, it is possible to focus the attention on the most informative sensor for the different time window of the various home. A possible application of the proposed methodology is to consider a reasonable period of time, i.e., 2 weeks, such to generate a behavioural model specific for each house. Each model can be used as comparison tool to identify if there any changes in the following daily activities of elderly representing a possible physical or mental illness that affect his or her health. In case that a change is detected, this information can be provided to a caregiver or family members that can assess the health of the elderly. The period for the creation of the behavioural model is going to be updated regularly such to consider the different season and the on-going aging of the inhabitants. Therefore, with this research, it is possible to state that the well-being of elderlies can be monitored, improved, and predicted using traditional methods and innovative approaches derived from the computer sciences. Future works regarding the smartwatch as assistive technology could be focused on a larger population of elderlies, including a heterogeneous age range or subjects affected by a progressive decline in physical and cognitive skills but still independent. A further step could be given by the combination of the results obtained by the analysis of domotic data with the e-coaching solutions to create a personalized smart home configuration with improved assistive and pro-active capabilities of improving the health of the inhabitants before that any decline of physical and mental could affect the health of the users.

# Bibliography

1. *Dynamics of life expectancy and life span equality*. **José Manuel Aburtoa, Francisco Villavicencio, Ugofilippo Basellini, Søren Kjærgaard, James W. Vaupel**. 2020, Proceedings of the National Academy of Sciences, Vol. 117.
2. **World Health Organization: Ageing and health**. <https://www.who.int/news-room/fact-sheets/detail/ageing-and-health>. [Online] 2018.
3. *World Population Prospects 2019*. **United Nations, Department of Economic and Social Affairs, Population Division**. 2019.
4. **World Health Organization. Disability and health**. <https://www.who.int/en/news-room/fact-sheets/detail/disability-and-health>. [Online] 2020.
5. *Economic and social implications of aging societies*. **Sarah Harper**. 2014, Science, Vol. 346.
6. *A Survey on Ambient-Assisted Living Tools for Older Adults*. **Parisa Rashidi, Alex Mihailidis**. 2013, IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS, Vol. 17.
7. *Long-Term Care for Older Adults: The Role of the Family*. **Janzen, Wonita**. 2001, Journal of Gerontological Nursing, Vol. 27.
8. *Introducing knowledge in the process of supervised classification of activities of Daily Living in Health Smart Homes*. **Anthony Fleury, Norbert Noury, Michel Vacher**. 2010. The 12th IEEE International Conference on e-Health Networking, Applications and Services.
9. *Smart homes for elderly healthcare—recent advances and research challenges*. **Sumit Majumder, Emad Aghayi, Moein Noferesti, Hamidreza Memarzadeh-Tehran, Tapas Mondal, Zhibo Pang, and M Jamal Deen**. 2017, Sensors, Vol. 17.
10. *Systematic Literature Review of Smart Home Monitoring Technologies Based on IoT for the Elderly*. **Kholoud Maswadi, Norjihani Binti Abdul Ghani, Suraya Binti Hamid**. 2020, IEEE Access, Vol. 8.

11. *How can technology support ageing in place in healthy older adults? A systematic review.* **Aline Ollevier, Gabriel Aguiar, Marco Palomino, Ingeborg Sylvia Simpelaere.** 2020, Public Health Reviews, Vol. 41.
12. *A Review of Internet of Things Technologies for Ambient Assisted Living Environments.* **Rytis Maskeliunas, Robertas Damaševičius, Sagiv Segal.** 2019, Future Internet, Vol. 11.
13. *Instrumentation and measurement in medical, biomedical, and healthcare systems.* **Shervin Shirmohammadi, Kurt Barbe, Domenico Grimaldi, Sergio Rapuano, Sabrina Grassini.** 2016, IEEE Instrumentation & Measurement Magazine, Vol. 19.
14. *An IoT Approach for an AAL Wi-Fi-Based Monitoring System.* **Marco Bassoli, Valentina Bianchi, Ilaria De Munari, Paolo Ciampolini.** 2017, IEEE Transactions on Instrumentation and Measurement , Vol. 66.
15. *Smartfaber: Recognizing fine-grained abnormal behaviors for early detection of mild cognitive impairment.* **Daniele Riboni, Claudio Bettini, Gabriele Civitarese, Zaffar Haider Janjua, Rim Helaoui.** 2016, Artificial Intelligence in Medicine, Vol. 67.
16. *Automated Cognitive Health Assessment Using Smart Home Monitoring of Complex Tasks.* **Prafulla N. Dawadi, Diane J. Cook, Maureen Schmitter-Edgecombe.** 2013, IEEE transactions on systems, man, and cybernetics: systems, Vol. 17.
17. *Ambient Sensors for Elderly Care and Independent Living: A Survey.* **Md. Zia Uddin, Weria Khaksar, Jim Torresen.** 2018, Sensors, Vol. 18.
18. *Exploring the ambient assisted living domain: A systematic review.* **Davide Calvaresi, Daniel Cesarini, Paolo Sernani, Mauro Marinoni, Aldo Franco Dragoni, Arnon Sturm.** 2017, Journal of Ambient Intelligence and Humanized Computing, Vol. 8.
19. *A Review of Ambient Intelligence Assisted Healthcare Monitoring.* **Abdelhamid Salih Mohamed Salih, Ajith Abraham.** 2013, International Journal of Computer Information Systems and Industrial Management Applications., Vol. 5.
20. *A Survey on Ambient Intelligence in Healthcare.* **Giovanni Acampora, Diane J. Cook, Parisa Rashidi, Athanasios V. Vasilakos.** 2013, Proceedings of the IEEE, Vol. 101.

21. *Detecting Health and Behavior Change by Analyzing Smart Home Sensor Data*. **Gina Sprint, Diane J. Cook, Roschelle Fritz, Maureen Schmitter-Edgecombe**. 2016. IEEE International Conference on Smart Computing (SMARTCOMP).
22. *Smart Homes for Older People: Positive Aging in a Digital World*. **Quynh Lê, Hoang Boi Nguyen, Tony Barnett**. 2012, Future Internet, Vol. 4.
23. **Joost van Hoof, George Demiris, Eveline J.M. Wouters**. *Handbook of Smart Homes, Health Care and Well-Being*. s.l. : Springer, 2017.
24. *Assistive sensor-based technology driven self-management for building resilience among people with early stage cognitive impairment*. **Sara Casaccia, Roberta Bevilacqua, Lorenzo Scalise, Gian Marco Revel, Arlene J.Astell, Susanna Spinsante, Lorena Rossi**. 2019, IEEE International Symposium on Measurements & Networking (M&N).
25. <http://www.aal-europe.eu>. [Online]
26. <https://resilien-t.eu/>. [Online]
27. **Sazonov, Edward**. *Wearable Sensors: Fundamentals, Implementation and Applications*. s.l. : Elsevier, 2020.
28. *Transformation in Healthcare by Wearable Devices for Diagnostics and Guidance of Treatment*. **Aman Mahajan, Gregory Pottie, William Kaiser**. 2020, ACM Transactions on Computing for Healthcare, Vol. 1.
29. *A Review of Wearable Technologies for Elderly Care that Can Accurately Track Indoor Position, Recognize Physical Activities and Monitor Vital Signs in Real Time*. **Zhihua Wang, Zhaochu Yang**. 2017, Sensors, Vol. 17.
30. *Wearable Health Devices—Vital Sign Monitoring, Systems and Technologies*. **Duarte Dias, João Paulo Silva Cunha**. 2018, Sensors, Vol. 18.
31. <https://www.sensara.care/index.php>. [Online]
32. <http://www.aal-europe.eu/projects/rosetta/>. [Online]



33. *Unobtrusive Health Monitoring in Private Spaces: The Smart Home*. **Ju Wang, Nicolai Spicher, Joana M. Warnecke, Mostafa Haghi, Jonas Schwartz, Thomas M. Deserno**. 2021, *Sensors*, Vol. 21.
34. *Healthcare wearable devices: an analysis of key factors for continuous use intention*. **Sang M. Lee, DonHee Lee**. 2020, *Service Business*, Vol. 14.
35. *A Survey on Wearable Technology: History, State-of-the-Art and Current Challenges*. **Aleksandr Ometov et al.** 2021, *Computer Networks*, Vol. 193.
36. **Samiya Khan, Mansaf Alam**. *Wearable Internet of Things for Personalized Healthcare: Study of Trends and Latent Research*. [book auth.] Anupam Biswas, Pinki Roy Ripon Patgiri. *Health Informatics: A Computational Perspective in Healthcare*. s.l. : Springer, 2021.
37. *Personalized Health Monitoring System of Elderly Wellness at the Community Level in Hong Kong*. **Lisha Yu, Wai Man Chan, Yang Zhao, Kwok-Leung Tsui**. 2018, *IEEE Access*, Vol. 6.
38. *Forecasting one-day-forward wellness conditions for community-dwelling elderly with single lead short electrocardiogram signals*. **Xiaomao Fan, Yang Zhao, Hailiang Wang, Kwok Leung Tsui**. 2019, *BMC Medical Informatics and Decision Making*, Vol. 19.
39. **Drijfhout Desmond**. *Predicting well-being with wearable sensor data*. University of Groningen. 2019.
40. *Monitoring and Prediction of Mood in Elderly People during Daily Life Activities*. **Daniel Bautista-Salinas, Joaquín Roca González, Inmaculada Méndez, Oscar Martínez Mozos**. 2019. Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC).
41. *Context Aware Computing for The Internet of Things: A Survey*. **Charith Perera, Arkady Zaslavsky, Peter Christen, Dimitrios Georgakopoulos**. 2014, *IEEE Communications Surveys & Tutorials*, Vol. 16.
42. *Internet of Things: A Survey on Enabling Technologies, Protocols, and Applications*. **Ala Al-Fuqaha, Mohsen Guizani, Mehdi Mohammadi, Mohammed Aledhari, Moussa Ayyash**. 2015, *IEEE Communications Surveys & Tutorials*, Vol. 17.

43. **Richard Harper** . *Inside the Smart Home: Ideas, Possibilities and Methods*. s.l. : Springer, 2003.
44. *Unsupervised Machine Learning for Developing Personalised Behaviour Models Using Activity Data*. **Laura Fiorini, Filippo Cavallo, Paolo Dario, Alexandra Eavis, Praminda Caleb-Solly**. 2017, *Sensors*, Vol. 17.
45. *Tracking changes in user activity from unlabelled smart home sensor data using unsupervised learning methods*. **Prankit Gupta, Richard McClatchey, Praminda Caleb-Solly**. 2020, *Neural Computing and Applications*, Vol. 32.
46. *Smart homes for the elderly dementia sufferers: identification and prediction of abnormal behaviour*. **Ahmad Lotfi, Caroline Langensiepen, Sawsan M. Mahmoud, M. J. Akhlaghinia**. 2012, *Journal of Ambient Intelligence and Humanized Computing*, Vol. 3.
47. *Unobtrusive Anomaly Detection in Presence of Elderly in a Smart-home Environment*. **Marek Novàk, Miroslav Binas, Frantisek Jakab**. 2012, *ELEKTRO*.
48. *Virtual Coaches for Older Adults' Wellbeing: A Systematic Review*. **Mira El Kamali et al.** 2020, *IEEE Access*, Vol. 8.
49. **www.ihealthlabs.it**. [Online]
50. **<https://ihealthlabs.eu/it/48-ihealth-wave-855111003965.html>**. [Online]
51. *A Novel Virtual Coaching System Based on Personalized Clinical Pathways for Rehabilitation of Older Adults—Requirements and Implementation Plan of the vCare Project*. **Sofoklis Kyriazakos et al.** 2020, *Frontiers in Digital Health*, Vol. 2.
52. *Smartwatches selection: market analysis and metrological characterization on the measurement of number of steps*. **Casaccia Sara, Revel Gian Marco, Scalise Lorenzo, Cucchieri Giacomo and Rossi Lorena**. 2021 , *IEEE MeMeA*.
53. *Listening to respondents: a qualitative assessment of the Short-Form 36 Health Status Questionnaire*. **Sara Mallison**. 2002, *Social Science & Medicine*, Vol. 54.
54. *The SF-36 health survey questionnaire: is it suitable for use with older adults?* **V. Hayes, J. Morris, C. Wolfe, M. Morgan**. 1995, *Age Ageing*, Vol. 24.

55. *Testing the measurement properties of the Short Form-36 Health Survey in a frail elderly population.* **K. Stadnyk, J. Calder, K. Rockwood.** 1998, Journal of Clinical Epidemiology, Vol. 51.
56. *Support-Vector Networks.* **Corinna Cortes, Vladimir Vapnik.** 1995, Machine Learning, Vol. 20.
57. *Scikit-learn: Machine Learning in Python.* **Fabian Pedregosa et al.** 2011, Journal of Machine Learning Research, Vol. 12.
58. **Jerome H. Friedman, Robert Tibshirani, Trevor Hastie.** *The elements of statistical learning: data mining, inference, and prediction.* 2. s.l. : Springer, 2009.
59. *Kernel Methods in Machine Learning.* **Thomas Hofmann, Bernhard Schölkopf, Alexander J. Smola.** 2008, The Annals of Statistics, Vol. 36.
60. **Shai Ben-David, Shai Shalev-Shwartz.** *Understanding Machine Learning: From Theory to Algorithms.* s.l. : Cambridge University Press, 2014.
61. *Top 10 algorithms in data mining.* **Xindong Wu et al.** 2008, Knowledge and Information Systems, Vol. 14.
62. *Top-down induction of decision trees classifiers - a survey.* **Lior Rokach, Oded Maimon.** 2005, IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS, Vol. 35.
63. *Random Decision Forests.* **Ho, Tin Kam.** 1995, International Conference on Document Analysis and Recognition.
64. *The random subspace method for constructing decision forests.* **Ho, Tin Kam.** 1998, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 20.
65. *Random Forests.* **Leo Breiman.** 2001, Machine Learning, Vol. 45.
66. *Learning Internal Representations by Error Propagation.* **Rumelhart David E., Hinton Geoffrey E., Williams Ronald J.** 1985, Parallel distributed processing: explorations in the microstructure of cognition, Vol. 1.
67. *Adam: A Method for Stochastic Optimization.* **Diederik Kingma, Jimmy Ba.** 2014.

68. *Evaluation: From Precision, Recall and F-Factor to ROC, Informedness, Markedness & Correlation*. **David M. W. Powers**. 2011, International Journal of Machine Learning Technologies.
69. *Selecting and interpreting measures of thematic classification accuracy*. **Stephen V. Stehman**. 1997, Remote Sensing of Environment, Vol. 62.
70. **Michael W. Berry, Azlinah Mohamed, Bee Wah Yap**. *Supervised and Unsupervised Learning for Data Science*. s.l. : Springer, 2020.
71. **Saurav Kaushick**. An Introduction to Clustering and different methods of clustering. [Online] <https://www.analyticsvidhya.com/blog/2016/11/an-introduction-to-clustering-and-different-methods-of-clustering/>.
72. **Berry Michael W., Azlinah Mohamed, Bee Wah Yap**. *Supervised and Unsupervised Learning for Data Science*. s.l. : Springer, 2020.
73. **M. Emre Celebi, Kemal Aydin**. *Unsupervised Learning Algorithms*. s.l. : Springer, 2016.
74. *Some Methods for classification and Analysis of Multivariate Observations*. **J. B. MacQueen**. 1967, Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability.
75. **Sanatan Mishra**. Unsupervised Learning and Data Clustering. [Online] <https://towardsdatascience.com/unsupervised-learning-and-data-clustering-eeecb78b422a>.
76. [https://matteucci.faculty.polimi.it/Clustering/tutorial\\_html/kmeans.html](https://matteucci.faculty.polimi.it/Clustering/tutorial_html/kmeans.html). [Online]
77. *Principal component analysis: a review and recent developments*. **Ian T. Jolliffe, Jorge Cadima**. 2016, Philosophical transactions of the royal society A, Vol. 374.
78. **Dubey Akash**. The Mathematics Behind Principal Component Analysis: From raw data to principal components. <https://towardsdatascience.com/the-mathematics-behind-principal-component-analysis-fff2d7f4b643>. [Online]
79. *Silhouettes: A graphical aid to the interpretation and validation of cluster analysis*. **Peter J. Rousseeuw**. 1987, Journal of Computational and Applied Mathematics, Vol. 20.

80. **Waisakurnia Winson.** The Easiest Way to Interpret Clustering Result.  
*<https://towardsdatascience.com/the-easiest-way-to-interpret-clustering-result-8137e488a127>*. [Online]

81. *Interpretable Clustering via Optimal Trees.* **Dimitris Bertsimas, Agni Orfanoudaki, Holly Wiberg.** 2018, arXiv e-prints.