

# Modelling Similarity for Comparing Physical Activity Profiles - A Data-driven Approach

Deepika Verma<sup>1</sup>, Kerstin Bach<sup>1</sup>, and Paul Jarle Mork<sup>2</sup>

<sup>1</sup> Department of Computer Science

<sup>2</sup> Department of Public Health and Nursing

Norwegian University of Science and Technology, Trondheim, Norway

<http://www.idi.ntnu.no>, <http://www.ntnu.no/ism>

**Abstract.** Objective measurements of physical behaviour are an interesting research field from the public health and computer science perspective. While for public health research, measurements with a high quality and feasible setup is important, the analysis of and reasoning about the data is what we will present in this work. Our focus in this work is the comprehensive representation of physical behaviour throughout consecutive days and allowing to find subgroups in the population with similar physical activity levels.

We have a unique data set of 4628 participants wearing tri-axial accelerometers for six days and will present a case-based reasoning (CBR) system that can find and compare similar activity profiles. In this work, we focus on creating a CBR model using myCBR and do initial experiments with the resulting system. We will introduce a data-driven approach for modelling local similarity measures. Eventually, in the experiments we will show that for the given data set, the CBR system outperforms a k-Nearest Neighbor regressor in finding most similar participants.

**Keywords:** Physical Activity, Case-Based Reasoning, Local Similarity Modelling, k-Nearest Neighbor

## 1 Introduction

Physical inactivity and poor sleep are considered global health problems [16,25] that contribute substantially to poor health and premature mortality. It is estimated that physical inactivity is responsible for about 9% premature mortality [19], which is similar to the effect of smoking [31] and obesity [1].

CBR has become more popular over the last few years, especially in an area where continuous measurements become more and more available [9,23]. It offers a way for abstracting and transferring specific domain expert knowledge into a self-explanatory and user-friendly tool, which can be used to generate solutions for problems ranging from simple daily life tasks to complex issues (which otherwise necessitate expert help), with an appropriate reasoning behind them. Not only is it being applied for finding similar cases to provide solutions, but also

for the classification of medical [33,8] and activity data [30]. In [30], the authors propose a CBR method to classify different physical activities of elderly based on their pulse rate.

In this paper, we focus on the knowledge engineering process of creating a CBR model and present a data-driven approach for modelling local similarity measures for physical activity data in the myCBR workbench [5,29]. We will show in our experiments that a CBR system comparing physical activity profiles is less erroneous than a k-Nearest Neighbour (k-NN) regressor model. In our experiments, both approaches are used to find groups of similar activity profiles and their performance is evaluated statistically. The second contribution of this paper is a method for modelling the local similarity measures utilizing data driven methods. We will showcase how a data set can lead to strong initial definitions for numerical value ranges and therewith ease and stratify the knowledge modelling process.

The remaining of this paper is divided into sections as follows: in section 2, we discuss related work on reasoning about physical activity behaviour using various approaches within machine learning and artificial intelligence. In section 3, we discuss the importance of objective measurements of physical activity behaviour from both public health and computer science perspective. Section 4 is dedicated to similarity modelling for the data set in myCBR. In section 5, we present the experiments performed to evaluate the CBR model generated and compare it with that of k-NN model. Section 6 and 7 are for discussion and conclusion respectively.

## 2 Related Work

The amalgamation of sensors, Internet of Things (IoT) and Artificial Intelligence (AI) provides a unique opportunity not only for health researchers, but also for AI researchers to perform objective measurements and utilize raw data recordings to generate physical activity profiles of a large number of participants and determine similar physical activity profile groups. With the help of AI techniques, it is possible to perform objective analysis of sensor data stream to not only identify different physical activities uniquely [7,4,32], but also find out groups of similar activity profiles. Finding and clustering similar physical activity profiles is crucial in facilitating the understanding of health and activity characteristics of a population and identifying different activity phenotypes<sup>3</sup>. In [21], the author proposed an ATLAS index to cluster and identify four activity phenotypes using NHANES<sup>4</sup> data set. Similarly, in [32], authors proposed a statistical machine learning model to identify different sleep and physical activity phenotypes. Further, the authors in [13] apply latent class analysis to identify five different activity phenotypes among young adults in a cohort study where data was collected using hip-worn accelerometers for seven days. Our long term

---

<sup>3</sup> <https://www.biology-online.org/dictionary/Phenotype>

<sup>4</sup> <https://wwwn.cdc.gov/nchs/nhanes/default.aspx>

goals and target data are similar to these studies, however the approach differs slightly.

Similar to the preference-based CBR framework presented by Hüllermeier and Schlegel [14], we are presenting a framework for modelling local similarity measures based on the data set available. Therewith we can tailor each similarity measure to the application domain. In the continuation of their work Abdel-Aziz, Strickert and Hüllermeier [2] show that the data distributions and distances in data sets can be used for learning similarity measures. While the authors focus on learning preferences, we show with the work presented here that the data-driven view can be carried over to general knowledge engineering tasks. Using a data-driven approach for automatic similarity learning and feature weighting has been presented by Gabel and Godehardt [11]. In their work they trained a neural network to induce local and global similarity measures. While we are not automatically assigning the similarity measures, we also use existing cases to derive them. In [28], the authors explore a case-based approach for recommending 5km times for marathon runners in order to achieve their personal best. The approach they apply is similar to the one presented in this paper as they use timing profiles as basis for the similarity-based assessment. In a slightly different approach, Sani et. al. [27] explore using k-NN for detecting physical activities from wrist worn sensors. In their work they show that applying k-NN for detecting movement patterns is very successful for creating personalized models. Even though the approaches differ, our work is similar in terms of comparing physical activity profiles with raw data coming from accelerometers.

### **3 Physical Activity Analysis for Public Health Application Scenarios**

Regular physical activity is important for people of all age groups, including the elderly. It can significantly reduce the risk of various health problems such as stroke, diabetes, various types of cancer, depression, as well as hypertension and improve bone and muscle health<sup>5</sup>. Physical inactivity is one of the most important public health problems of this century and has a strong negative impact on the physical and mental well being of an individual. It is estimated that about 23% adults and 81% adolescents globally are physically inactive. The figures are alarmingly high for adolescents. Moreover, being physically active is not just about moving around in the house or walking at a slow pace, they must include some form of Moderate to Vigorous Physical Activity (MVPA) such as brisk walking, dancing, running, cycling, or moving/lifting heavy load.

Over the last few years, researchers in public health domain have moved rapidly from using self-reported subjective activity data to objectively measured activity data with the use of body-worn sensors [4,18,20]. Not only are the sensors a more viable option due to the simplicity of extracting and utilizing raw data, but also eliminate the problem of bias due to self reporting [24,17], which

---

<sup>5</sup> [http://who.int/features/factfiles/physical\\_activity/en/](http://who.int/features/factfiles/physical_activity/en/)

has been a major concern among researchers as it leads to inaccuracy and uncertainty. Moreover, the accelerometers directly measure the subject's physiology motion status to indicate the motion pattern within a given time period, which is helpful in activity recognition and are much more energy efficient.

The physical activity data used for this work is primarily based on accelerometer data collected during the HUNT4<sup>6</sup> cohort study. The N rd-Tr ndelag Health Study (HUNT)<sup>7</sup> in Norway is one of the largest health studies of its kind. The study consists of a large amount of health data collected through questionnaires and clinical examinations during three intensive previous studies (HUNT1 1984-86, HUNT2 1995-97 and HUNT3 2006-08). In the ongoing study HUNT4 (2017-19), each participant is offered to participate in the objective measurements data collection. If accepted, they are fitted with two wearable tri-axial accelerometers, placed at their thigh and lower back, which are used to collect activity data for one week. The raw sensor data is then classified into 17 different physical activities using Support Vector Machines (for the synchronization of sensor data) and Random forest classifiers (for the prediction of activity classes). Afterwards, these activities are grouped into six main physical activities: lying, sitting, standing, walking, running, cycling, which is the basis data set for our work <sup>8</sup>.

By determining the variation among participants in different activity clusters through similarity, it is possible to provide activity recommendations to less active profiles in order to make them more active. Every person has different activity characteristics and finding a group of activity profiles most similar to that person with respect to the duration of every activity is a challenging task and we aim to address this task using Case-Based Reasoning (CBR), because it offers the flexibility and transparency in its reasoning process.

## 4 Data-driven Knowledge Modelling

In this section, we explain how we implement a CBR system that can be applied to find and compare similar activity profiles from objectively measured population data. We are using the local-global-principle [26] for creating similarity measures and thereby build a knowledge model that tailors the similarity measure for each attribute. Once the local similarity measures are defined, we continue to use weighted sum for defining the global similarity.

While the HUNT4 data set is unique in the world, the challenges for utilizing it for developing a CBR system are very common such as the identification of suitable data set context for the problem at hand, definition of initial similarity measures, representation of cases and determination of valuable cases for populating the casebase. In this work we will introduce a method for utilizing a given data set to model similarity measures. Further we will take into account the effect of growing casebases and show a methodology that can help to visualize and understand how a CBR system learns.

---

<sup>6</sup> <https://www.ntnu.no/hunt4/>

<sup>7</sup> <https://www.ntnu.no/hunt/>

<sup>8</sup> Since the study is ongoing, we have used the data available by March, 12 2018.

This section is further divided into subsections as follows: First, we describe how we populate the casebase and generate cases in the developed case representation. Second, we describe our data-driven approach to model the local similarity measures for the numerical activity attributes. Once the model is in place, we then query the casebase and compare the most similar activity profiles retrieved.

#### 4.1 Case Generation

Developing a case representation is the first part of the system development. Depending on the domain and the available data this can be a challenging process on its own [12,6,15]. For our application domain we utilize the pre-processed HUNT4 data. While HUNT4 collects a very comprehensive set of data, we are only focusing on the objective measurements. The sensor data is collected over a period of seven days per participant and the overall data collection in the cohort stretches over 18 months, starting from the autumn of 2017 until February 2019. It is an ongoing study and until March 2018, data for 17409 participants has been automatically classified and for each participant aggregated into the six main physical activities. In Table 1 we present the description of the six activity types used in our data set.

Activity	Description
Lying	The person lies down
Sitting	When the person's buttocks is on the seat of the chair, bed or floor
Standing	Upright, feet supporting the person's body weigh
Walking	Locomotion towards a destination with one stride or more
Running	Locomotion towards a destination, with at least two steps where both feet leave the ground during each stride
Cycling	The person is riding bicycle

**Table 1.** Activity Descriptions

Each participant is fitted with two tri-axial accelerometers, AX3 Axivity<sup>9</sup>, one on the thigh and second on lower back. The sensors are used to detect vibrations, movement and orientation changes in the three axes. The sampling frequency of the sensors is set at 50Hz. After the participant has worn the sensors for seven days, they are returned to the HUNT research center where the raw data is downloaded, extracted and classified using Support Vector Machines and Random Forest algorithms. The resulting data set contains the H4ID (unique ID for each HUNT4 participant), number of minutes of each different activity, the date and day of the week in a csv file.

When preparing the data for the CBR system, we further process it by removing the records where we assume the sensor was taken off or the prediction

<sup>9</sup> <https://axivity.com/downloads/ax3>

failed. Those are very long times of the same activities. Records are removed based on the following criteria:

- sum of all the activities for a single record exceeds 1440, which is the total minutes in a day
- records containing zero minutes for lying, sitting, standing and walking
- data set for one participant has less than seven days of data

Eventually, we chose to keep records where exactly six days of data per H4ID was present, while the rest of the records were removed. For each unique H4ID, the total minutes of each activity were summed up for six days. We experimented with different knowledge representations including mean, maximum and sum of duration of each activity per H4ID and found the sum representation to suit best since it captures the overall physical behaviour of the participants over the days as well as the variance of the similarity measure over its' entire range. At this point, after pre-processing, the data set contains 4628 rows, each record containing sum of each activity over six days for a single participant. Table 2 gives a brief account of the data set.

	Lying	Sitting	Standing	Walking	Running	Cycling
count	4628	4628	4628	4628	4628	4628
mean	3090.49	3322.82	1401.22	790.67	6.86	26.45
min	7.35	253.25	56.50	1.55	0	0
max	7513.80	7846.10	4247.10	2101.65	172.70	719.10

**Table 2.** Data set Statistics

Cases are populated from the previously described data set by loading into the previously defined case representation using the myCBR tool. A single case in myCBR is represented as shown in Fig. 1, where *Participant* is the name of the concept which consists of six attributes namely *cycling*, *lying*, *running*, *sitting*, *standing* and *walking*.

## 4.2 Data-driven Similarity Measures Development

The local-global-principle requires that both types of similarity measures, the local one on the attribute level and the global one on the conceptual need to be defined.

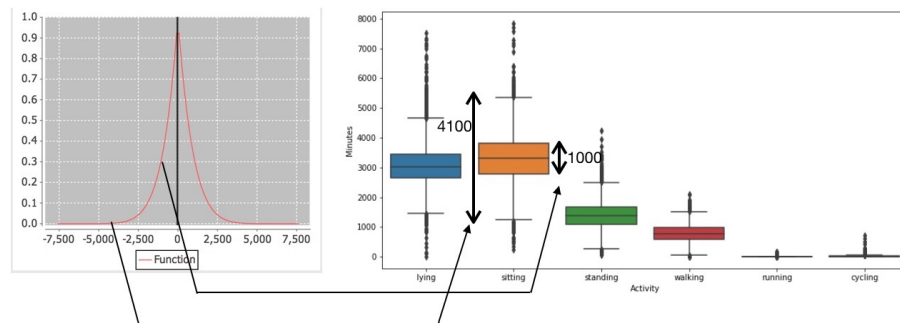
Modelling the local similarity measures for different attributes in myCBR can be challenging as researchers have to balance the input from the domain experts and the available data. Having criteria which can lead the knowledge modelling process is helpful for both parties. We therefore suggest to make use of the existing data in this process. As we assume that the collected data set covers the scope of what type of problems (cases) we have seen before, this is a useful departure point. In the following, we would have a reality check with

Instance information	
Name	Participant1
Attributes	
cycling	87.0
lying	3624.65
running	1.95
sitting	2819.35
standing	1258.75
walking	848.3

**Fig. 1.** Case representation in myCBR

the domain experts that discusses whether the defined value ranges cover the domain well. While setting upper and lower limits is straight forward, assigning the similarity behaviour is not. Consecutively, we assume that numerical local similarity measures are distance functions and the question is how steep of a similarity decline should be chosen. We use polynomial functions to model similarity measure since they are more flexible and provide better convergence when using continuous numerical data. Therefore, we will focus on the polynomial function of the similarity measure and our goal is to determine their degree.

Taken this task in our application domain, we see an activity variation among different profiles, but also in the aggregation of activities over all profiles. We use box plots for visualizing the distributions and variations in our data set and transfer this into modelling local similarity measures.



**Fig. 2.** Example for Data-driven Local Similarity Modelling: On the left there is a screen shot of a polynomial similarity function for a value range between 0 and 7500. With the arrows we depict how the box-plot for sitting relates to the decrease of similarity at a certain distance.  $IQR * 1.5$  method has been used for the box plots.

Fig 2 shows an example of a numerical local similarity measure. In the example, it is the total amount of sitting during six days. From there we look into the  $Q_1$  and  $Q_3$  which indicated the majority spread for the data set. We decided to take these values as reference points for determining the decrease of similarity.

Hence, creating a box-plot of the data set will allow modelling each activity attribute since we only take the Inter Quartile Range (IQR) and the range (min to max) into account:

$$\begin{aligned} r_1 &= IQR \\ r_2 &= range \end{aligned} \tag{1}$$

It represents the difference between upper ( $Q_3$ ) and lower ( $Q_1$ ) quartiles in the box-plot, that is  $IQR = Q_3 - Q_1$ .

We assume that all similarity functions are polynomial and adjust the polynomial degree of the similarity function such that

$$\begin{aligned} y(r_1) &\approx 0.30 \\ y(r_2) &\approx 0 \end{aligned} \tag{2}$$

We can observe in Fig 2 how the similarity function varies after applying the methodology in equation 1 and 2. The bigger the polynomial degree, the steeper the similarity function and more precise the attribute values in retrieved cases. The decline in the similarity function is steeper in the beginning until at  $r_1$  it reaches close to  $y(r_1)$  and then decreases gradually until at  $r_2$  it is approximately close to  $y(r_2)$ . This way, the similarity function covers the entire attribute range as well as the similarity measure range  $[0, 1]$ . While the choice of  $y(r_1)$  and  $y(r_2)$  depends on the domain-expert's knowledge and satisfaction with the outcome, we however experimented with different values and found these best suited for our application domain. We use this as the initial definition of similarity measures. If required, the function can of course be further customized if the relevant domain knowledge is available.

### 4.3 Comparing Physical Activity Profiles

Once the casebase and similarity measures are in place, the model can be used to find similar profiles. Fig 3 shows the result of one such query retrieval in myCBR. The figure shows that the retrieved cases are sorted by similarity value in descending order, that is, most similar case are displayed at the top while least similar are at the bottom. On the lower part of the screen shot the four most similar profiles are shown in a detailed view. The tool marks closer matches darker.

While the myCBR workbench indicates that we can do a similarity-based retrieval, it is hard to judge how the CBR system works with increasing casebase or changing similarity measures. In the next section we will investigate how the casebase size and different retrieval methods perform in our application domain.



Query		Special Value:	
cycling	17	none	Participant0 - 0.98
lying	4295	none	Participant1 - 0.91
running	20	none	Participant2680 - 0.8
sitting	2391	none	Participant2023 - 0.79
standing	1244	none	Participant208 - 0.79
walking	670	none	Participant3056 - 0.79
			Participant2470 - 0.78
			Participant3572 - 0.77
			Participant556 - 0.77
			Participant1302 - 0.76
			Participant1522 - 0.76
			Participant1228 - 0.76
			Participant2369 - 0.75

	Participant0	Participant1	Participant2680	Participant2023
Similarity	0.98	0.91	0.8	0.79
cycling	17.35	17.35	19.7	25.85
lying	4295.75	4295.0	4216.0	4421.15
running	20.85	20.85	9.0	13.9
sitting	2391.25	2391.0	2504.05	2363.25
standing	1244.1	1244.1	1210.85	1241.2
walking	670.7	848.3	680.4	574.65

Fig. 3. A Query and its retrieval result in the myCBR workbench

## 5 Evaluation of Increasing Casebase Sizes and Retrieval Methods

A performance evaluation of the CBR model has been conducted using holdout-repeat cross-validation in which 200 random cases were held out to be used for testing. Therewith for each run our casebase consisted of 4428 cases. A test set, comprising of ten randomly selected cases from the held out set of 200 cases, represents a single epoch in the experiments and performance is reported using Mean Relative Error (MRE) as a measure of precision. Each experiment is repeated five times and the results are averaged over all the epochs.

For each query instance  $q_i$  in the test set, the number of similar cases retrieved  $r$  from the casebase is 20. The relative error of each activity is the computed between  $q_i$  and  $r$  for one case at a time. The errors are averaged to obtain MRE of each activity for  $q_i$ . The process is repeated for every  $q_i$  in the test set, that is, for  $i = [1, 10]$ .

The MRE of the six activities are added to get the total relative error for each  $q_i$ . MRE is then calculated by averaging the relative errors for the entire queried test set.

The total relative error  $T$  for each queried instance is calculated as:

$$T = \sum_{i=1}^6 MRE(A_i)$$

where  $A$  is the activity type as they were introduced in section 4.1. MRE for the each test set is calculated as:

$$MRE = \frac{\sum_{i=1}^{10} T_i}{10}$$

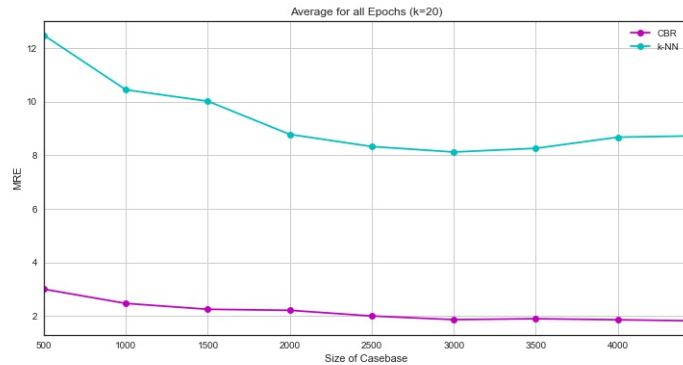
The experiments in this evaluation are performed in two ways: First, by calculating the MRE of retrieved instances against each queried test instance

with increasing casebase size. Second, by comparing the different results obtained using the CBR model and k-NN regressor model.

### 5.1 Increasing Casebase Size

This experiment focuses on the variation observed in MRE with the increasing size of the casebase. The CBR model was implemented using myCBR, however the tool does not support batch queries, which was the need of the hour for conducting the experiments for our work. To overcome this limitation, we used a myCBR Rest API <sup>10</sup> for batch querying the casebase using POST calls and the implementation was done in Python (version 3.6.3).

In this experiment, a test set is passed as a query using POST call when the casebase initially has 500 instances. Subsequently, MRE for that test set is calculated. 500 cases are then added to the casebase and the process is repeated until the casebase consists of the entire data set. The experiment is repeated five times, each with a different random test set. The average MRE of all the epochs for the given casebase size is shown in Fig 4.



**Fig. 4.** MRE comparison between the CBR model and k-NN regressor model with increasing casebase sizes (MRE is calculated for  $k = 20$  retrieved cases)

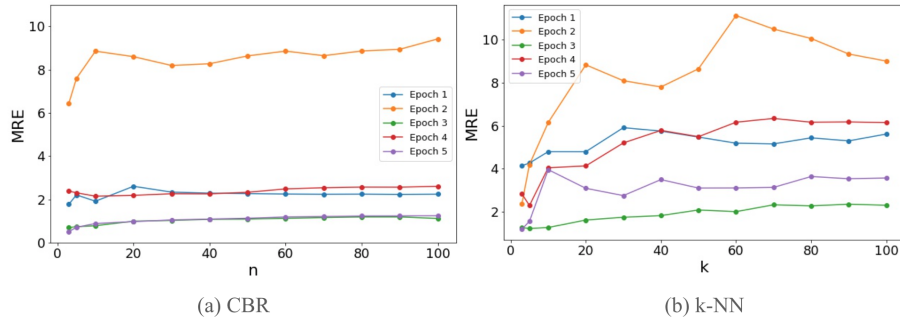
In order to have a comparison of the performance of the CBR model, the same experiment was conducted using k-NN regression model (with  $k = 20$ ). The implementation of the k-NN regressor was done using Scikit learn [22] library (version 0.19.1) in Python (version 3.6.3). The results obtained with the k-NN model are presented along with the results of the CBR model in Fig 4, where x-axis shows the size of the casebase (or size of data set for k-NN) and y-axis shows the MRE averaged over five epochs.

<sup>10</sup> <https://github.com/kerstinbach/mycbr-rest-example>

It can be observed from the results that MRE decreases steadily with increase in size of the casebase in the CBR implementation. However, the same cannot be said for k-NN, as the results show uncertain response to the increase in size of the data set. Even after introducing the entire data set, no improvement is observed. This decline in performance in k-NN is caused by the presence of outliers in the test set. CBR is able to estimate closest similar cases with respect to every activity for outliers very well, whereas k-NN cannot estimate the nearest neighbors with respect to every activity when presented with outliers. For instance, if there is an instance in the test set which has some or all attributes with values either below 25% or above 75% of the data range for those attributes in the data set, it leads to the k-NN algorithm computing nearest neighbors which are closer to the non-outlier attributes but farther from the outlier attributes. Thus, resulting in higher MRE even with an increased size of the data set.

## 5.2 Selection of k

Selecting an appropriate value of  $k$  is crucial in determining the success or failure of a k-NN regressor model. To see how the error varies, we experimented with different values of  $k$  in the range [3,100]. Fig 5(b) shows the variation in MRE with the change in value of  $k$ . Here, x-axis shows the value of  $k$  and y-axis shows the MRE.



**Fig. 5.** Number of closest cases: On the left is the graph depicting the variation in MRE with the number of most similar cases retrieved ( $n$ ) in CBR implementation. On the right is the graph for k-NN model depicting the variation in MRE with different values of  $k$ .

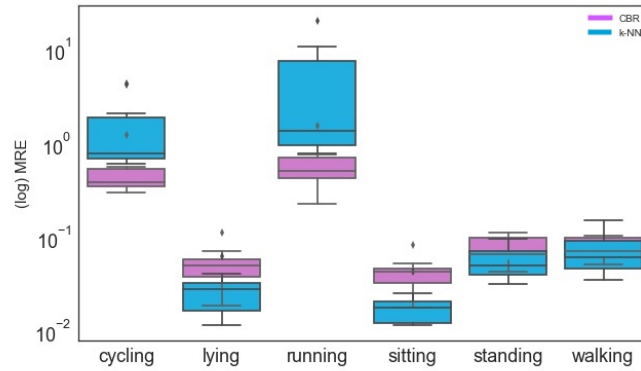
Although the determination of the closest similar profile in the CBR model is independent of  $n$  (number of retrieved cases), it is interesting to see how the MRE changes by varying  $n$  progressively. This allows us to further compare and contrast the performance of CBR model with that of k-NN model. Fig 5(a) shows the variation of MRE with increasing value of  $n$  in myCBR, where the x-axis shows the value of  $n$  and y-axis shows the MRE. It is clear from the

results that the value of  $k$  in k-NN (refer Fig 5(b)) has a huge impact on the MRE for each epoch. The implication of this graph is that with an increase in  $k$ , more neighboring cases are taken into consideration which are either less similar altogether or less similar with respect to a subset of activities, resulting in the sudden variation in errors. Whereas the CBR model has a relatively smoother response in creating the number of retrieved similar cases. It can be argued from the results that lower values of  $k$  would have been more suitable due to less MRE. However, our aim in this work is not to predict using k-NN, but to find a number of nearest neighbors of the queried profile, which is why we chose  $k = 20$  for our experiments. As our data set is large,  $k = 20$  is reasonably acceptable for this application domain. Also, from CBR perspective, considering more neighboring profiles helps in making improvements to the similarity measure to a greater extent than considering just one neighbor profile.

### 5.3 Composition of Error

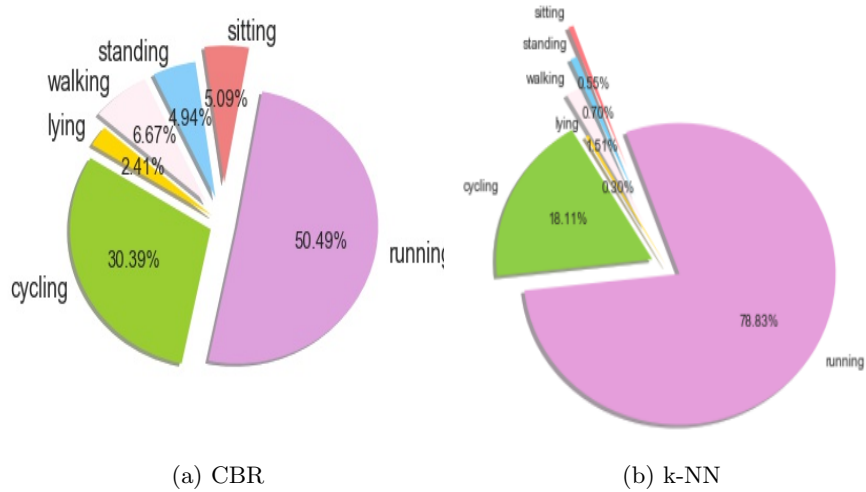
As we are using activity data to find other similar profiles, it is important to know the error observed in the approximation of each activity in the similar profiles.

Fig 6 shows the MRE (in log) for each activity using both the approaches when introduced with the entire data set. The figure underlines that for inactive time (lying, sitting, standing) - which is the majority for the participants (see Table 2 and Fig 2) - the k-NN approach produces less of an error. For moderate activities, like walking, both approaches are very close, while for rigorous activities, which we see only limited in the data set, the CBR approach produces much better results. This is very important for our overall aim of this work, as we eventually want to identify beneficial physical activity phenotypes.



**Fig. 6.** MRE per activity for the entire data set by the k-NN regressor and the CBR model

This observation is undermined by Fig 7, which shows the distribution of MRE for each of the activity calculated for both approaches after introducing the entire data set. It can be observed that in both k-NN and CBR, most of the error is attributed to the approximation of activity *running* (approx. 79% and 51% respectively). On the other hand, it is far lower in CBR, the result of which is relatively higher error composition of other activities as compared to those in k-NN. However, since these are compositional parts and convey only relative information, rather than concrete information, we must take into consideration the actual MRE, refer to Fig 4, which is significantly lower in case of CBR.



**Fig. 7.** Error Composition for the CBR (a) and k-NN (b) model

## 6 Discussion

The experimental results shown in Fig 4 demonstrate that the CBR model performs well in finding similar physical activity profiles. While k-NN is able to well approximate four out of six physical activities when finding the nearest neighbours, however it fails miserably in finding with respect to the other two activities, which results in higher MRE. On the other hand, the CBR model is able to determine the most similar physical activity profiles with respect to every activity more closely, resulting in far lower MRE as compared to the k-NN model. Furthermore, k-NN is susceptible to outliers, which is the cause of increase in MRE even after introducing the entire data set. Whereas this is not an issue with the CBR model. In Fig 5 we observe very minor increase in MRE with increasing number of retrieved instances using CBR model, whereas the

variations are more pronounced when using the k-NN model. These experiments demonstrate that the similarity modelling approach presented is working successfully for our application domain. Consequently, the CBR model significantly outperforms the k-NN algorithm and is more robust in finding similar physical activity profiles in a population. CBR approach can be applied to find and cluster similar activity groups, which will further be helpful in determining activity phenotypes.

## 7 Conclusion and Future Work

In this paper, we presented an approach to model the local similarity measures for physical behaviour data in myCBR in a data-driven manner. This model can be applied on physical behaviour data acquired using wearable sensors to find, group and compare similar activity profiles. We have demonstrated through experiments and statistical evaluation how the CBR model outperforms the state-of-the-art k-NN regressor model. Thus, it can be concluded that CBR approach is a suitable and viable option for application such as this in the public health domain. It can further be utilized in determining activity phenotypes in order to provide personalized activity recommendations to participants and help slowly transform an inactive into a more active lifestyle. We have also demonstrated through experiments the effectiveness of similarity modelling approach presented in this paper for the public health domain and it will be safe to conclude that it can be transferred to other similar domains dealing with continuous numerical data.

The method presented can further be enhanced to automatically assign the local similarities based on the attributes' values in the casebase using machine learning techniques, similar to what [11] presented in their paper. It can significantly reduce the efforts required to create new CBR models using different data sets from scratch.

In the future, we aim to extend our research towards compositional data analysis [3] on the HUNT4 data and applying CBR on the resulting compositional data. Compositional data analysis has been applied by researchers [10] for estimating the effect of change in physical activity behaviour for daily activities. Whether a change in one type of behaviour is beneficial or harmful for health depends on the compensatory shifts in other behaviours. The compositional nature of the HUNT4 data has therefore important consequences for both the analytical approach undertaken and interpretation of effects on health outcomes. Utilizing CBR for compositional data analysis will facilitate (i) getting insights into the behavioural characteristics between similar profiles in a population, (ii) understanding the association and co-dependency among various behaviours in different profiles, and (iii) identifying physical behaviour phenotypes.

## References

1. A, A., MH, F., MB, R.: Health effects of overweight and obesity in 195 countries over 25 years. *New England Journal of Medicine* 377(1), 13–27 (2017), pMID:

2. Abdel-Aziz, A., Strickert, M., Hüllermeier, E.: Learning solution similarity in preference-based cbr. In: Lamontagne, L., Plaza, E. (eds.) *Case-Based Reasoning Research and Development*. pp. 17–31. Springer International Publishing, Cham (2014)
3. Aitchison, J., Egozcue, J.J.: Compositional data analysis: Where are we and where should we be heading? *Mathematical Geology* 37(7), 829–850 (2005)
4. Arif, M., Kattan, A.: Physical activities monitoring using wearable acceleration sensors attached to the body. *PLOS ONE* 10(7), 1–16 (2015)
5. Bach, K., Althoff, K.D.: Developing Case-Based Reasoning Applications Using myCBR 3. In: Watson, I., Agudo, B.D. (eds.) *Case-based Reasoning in Research and Development, Proceedings of the 20th International Conference on Case-Based Reasoning (ICCBR-12)*. pp. 17–31. LNAI 6880, Springer (2012)
6. Bergmann, R., Kolodner, J., Plaza, E.: Representation in case-based reasoning. *The Knowledge Engineering Review* 20(03), 209 (2005)
7. Bulling, A., Blanke, U., Schiele, B.: A tutorial on human activity recognition using body-worn inertial sensors. *ACM Computing Surveys* 46(3), 1–33 (2014)
8. Campillo-Gimenez, B., Jouini, W., Bayat, S., Cuggia, M.: Improving case-based reasoning systems by combining k-nearest neighbour algorithm with logistic regression in the prediction of patients’ registration on the renal transplant waiting list. *PLoS ONE* 8(9) (2013)
9. Canensi, L., Leonardi, G., Montani, S., Terenziani, P.: Multi-level interactive medical process mining. In: ten Teije, A., Popow, C., Holmes, J.H., Sacchi, L. (eds.) *Artificial Intelligence in Medicine*. pp. 256–260. Springer, Cham (2017)
10. Dumuid, D., Pedisic, Z., Stanford, T.E., Martín-Fernández, J.A., Hron, K., Maher, C.A., Lewis, L.K., Olds, T.: The compositional isotemporal substitution model: A method for estimating changes in a health outcome for reallocation of time between sleep, physical activity and sedentary behaviour. *Statistical Methods in Medical Research* (2017)
11. Gabel, T., Godehardt, E.: Top-down induction of similarity measures using similarity clouds. In: Hüllermeier, E., Minor, M. (eds.) *Case-Based Reasoning Research and Development*. pp. 149–164. Springer International Publishing, Cham (2015)
12. H. El-Sappagh, S., Elmogy, M.: Case representation and indexing. *Foundations of Soft Case-Based Reasoning* p. 34–74 (2004)
13. Howie, E.K., Smith, A.L., Mcveigh, J.A., Straker, L.M.: Accelerometer-derived activity phenotypes in young adults: a latent class analysis. *International Journal of Behavioral Medicine* (2018)
14. Hüllermeier, E., Schlegel, P.: Preference-based cbr: First steps toward a methodological framework. In: Ram, A., Wiratunga, N. (eds.) *Case-Based Reasoning Research and Development*. pp. 77–91. Springer, Berlin, Heidelberg (2011)
15. Khamparia, A., Pandey, B.: A novel method of case representation and retrieval in cbr for e-learning. *Education and Information Technologies* 22(1), 337–354 (2017)
16. Kohl, H.W., Craig, C.L., Lambert, E.V., Inoue, S., Alkandari, J.R., Leetongin, G., Kahlmeier, S.: The pandemic of physical inactivity: global action for public health. *The Lancet* 380(9838), 294–305 (2012)
17. Lagersted-Olsen, J., Korshøj, M., Skotte, J., Carneiro, I., Søggaard, K., Holtermann, A.: Comparison of objectively measured and self-reported time spent sitting. *International Journal of Sports Medicine* 35(06), 534–540 (2013)
18. Lee, I.M., Shiroma, E.J.: Using accelerometers to measure physical activity in large-scale epidemiological studies: issues and challenges. *British Journal of Sports Medicine* 48(3), 197–201 (2013)

19. Lee, I.M., Shiroma, E.J., Lobelo, F., Puska, P., Blair, S.N., Katzmarzyk, P.T.: Effect of physical inactivity on major non-communicable diseases worldwide: an analysis of burden of disease and life expectancy. *The Lancet* 380(9838), 219–229 (2012)
20. Li, X., Dunn, J., Salins, D., Zhou, G., Zhou, W., Rose, S.M.S.F., Perelman, D., Colbert, E., Runge, R., Rego, S., et al.: Digital health: Tracking physiomes and activity using wearable biosensors reveals useful health-related information. *PLOS Biology* 15(1) (2017)
21. Marschollek, M.: A semi-quantitative method to denote generic physical activity phenotypes from long-term accelerometer data – the atlas index. *PLoS ONE* 8(5) (2013)
22. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12, 2825–2830 (2011)
23. Plis, K., Bunesco, R.C., Marling, C.R., Shubrook, J., Schwartz, F.: A machine learning approach to predicting blood glucose levels for diabetes management. In: *AAAI Workshop: Modern Artificial Intelligence for Health Analytics* (2014)
24. Prince, S.A., Adamo, K.B., Hamel, M., Hardt, J., Gorber, S.C., Tremblay, M.: A comparison of direct versus self-report measures for assessing physical activity in adults: a systematic review. *International Journal of Behavioral Nutrition and Physical Activity* 5(1), 56 (2008)
25. Raitakan, O.T., Porkka, K.V.K., Taimela, S., Telama, R., Räsänen, L., Vllkari, J.S.: Effects of persistent physical activity and inactivity on coronary risk factors in children and young adults the cardiovascular risk in young finns study. *American Journal of Epidemiology* 140(3), 195–205 (1994)
26. Richter, M.M.: The knowledge contained in similarity measures. In: Veloso, M.M., Aamodt, A. (eds.) *Case-Based Reasoning Research and Development*, Proc of the First International Conference, ICCBR-95. LNCS, vol. 1010. Springer (1995)
27. Sani, S., Wiratunga, N., Massie, S., Cooper, K.: knn sampling for personalised human activity recognition. In: Aha, D.W., Lieber, J. (eds.) *Case-Based Reasoning Research and Development*. pp. 330–344. Springer (2017)
28. Smyth, B., Cunningham, P.: Running with cases: A cbr approach to running your best marathon. In: Aha, D.W., Lieber, J. (eds.) *Case-Based Reasoning Research and Development*. pp. 360–374. Springer International Publishing, Cham (2017)
29. Stahl, A., Roth-Berghofer, T.R.: Rapid prototyping of cbr applications with the open source tool mycbr. In: *ECCBR '08: Proc. of the 9th European conference on Advances in Case-Based Reasoning*. pp. 615–629. Springer, Berlin (2008)
30. Uddin, M., Loutfi, A.: Physical activity identification using supervised machine learning and based on pulse rate. *International Journal of Advanced Computer Science and Applications* 4(7) (2013)
31. Wen, C.P., Wu, X.: Stressing harms of physical inactivity to promote exercise. *The Lancet* 380(9838), 192–193 (2012)
32. Willetts, M., Hollowell, S., Aslett, L., Holmes, C., Doherty, A.: Statistical machine learning of sleep and physical activity phenotypes from sensor data in 96,220 uk biobank participants. *BioRxiv* (2018)
33. Yao, B., Li, S.: Anmm4cbr: a case-based reasoning method for gene expression data classification. *Algorithms for Molecular Biology* 5(1), 14 (2010)