

Human Activity Recognition based on Wearable Sensor Data: A Benchmark

Artur Jordao, Antonio Carlos Nazare, Jessica Sena, William Robson Schwartz

Smart Sense Laboratory, Computer Science Department, Universidade Federal de Minas Gerais, Brazil

Email: {arturjordao, antonio.nazare, jessicasena, william}@dcc.ufmg.br

Abstract—Human activity recognition based on wearable sensor data has been an attractive research topic due to its application in areas such as healthcare and smart environments. In this context, many works have presented remarkable results using accelerometer, gyroscope and magnetometer data to represent the activities categories. However, current studies do not consider important issues that lead to skewed results, making it hard to assess the quality of sensor-based human activity recognition and preventing a direct comparison of previous works. These issues include the samples generation processes and the validation protocols used. We emphasize that in other research areas, such as image classification and object detection, these issues are already well-defined, which brings more efforts towards the application. Inspired by this, we conduct an extensive set of experiments that analyze different sample generation processes and validation protocols to indicate the vulnerable points in human activity recognition based on wearable sensor data. For this purpose, we implement and evaluate several top-performance methods, ranging from handcrafted-based approaches to convolutional neural networks. According to our study, most of the experimental evaluations that are currently employed are not adequate to perform the activity recognition in the context of wearable sensor data, in which the recognition accuracy drops considerably when compared to an appropriate evaluation approach. To the best of our knowledge, this is the first study that tackles essential issues that compromise the understanding of the performance in human activity recognition based on wearable sensor data.

Index Terms—Human activity recognition, wearable sensor data, state-of-the-art benchmark.

I. INTRODUCTION

Due to technological advances, activity recognition based on wearable sensors has attracted a large number of studies. This task consists of assigning a category of activity to the signal provided by wearable sensors such as accelerometer, gyroscope and magnetometer. With this purpose, there are two types of approaches. One investigates how to represent the raw signal to better distinguish the activities [1], [2] and the other explores the classification stage, which might consider either the raw signal [3], [4], [5], [6], a pre-defined representation [7], [8], [9], [10] or a combination of both [11], [12]. Even though many improvements have been achieved in both approaches [13], [14], [15], it is hard to assess the quality of sensor-based human activity recognition since previous studies have not contemplated the following issues.

- 1) Since many works propose their own datasets to conduct the evaluation, it is unclear which datasets are more adequate. In addition, these datasets are not always

publicly available, preventing the reproducibility of the work.

- 2) The metrics and the validation protocol employed to assess the activity recognition quality vary from paper to paper (e.g., the third and fourth columns in Table I).
- 3) The process to generate the data samples, before performing the evaluation, presents a wide variability (e.g., the second column in Table I).

While the first two issues play an important role, the last is a critical point, influencing the final performance of the activity recognition. It is important to emphasize that in other research areas, such as object detection [16], image classification [17] and face verification [18], [19], these issues are handled, leading to a standardized evaluation, which attracts more research efforts towards the application. As a consequence of these issues, currently it is not possible to know the state-of-the-art methods in human activity recognition based on wearable sensor data. Additionally, since the performance of the methods can be skewed by the conducted evaluation, it is difficult to identify promising research directions.

The aforementioned discussion motivated our study, in which we implement several existing works and evaluate them on a large set of publicly available datasets. To provide a more robust evaluation and show statistical equivalence between the methods, we also perform a statistical validation [20].

The development of this work presents the following contributions: 1) implementation and evaluation of several top-performance methods to human activity recognition based on wearable sensors, ranging from handcrafted-based to convolutional neural networks approaches; 2) demonstration that the process chosen to generate the data samples is a crucial point since all methods have their recognition accuracy reduced when evaluated on a unbiased sample generation process; 3) proposition of two novel data sample generation processes: *Full-Non-Overlapping-Window* and *Leave-One-Trial-Out*, in which each one handles a particular drawback of the currently employed sample generation processes; and 4) standardization of popular datasets focused on human activity recognition associated with wearable sensor data to facilitate their use and the evaluation of future works.

Since nowadays there are many works in human activity recognition based on wearable sensors, we selected the ones that provide enough information to reproducibility (i.e., the definition of the features employed and the classifier setup). In addition, regarding the datasets employed in our study, we select those used by previous works [10], [11], [8], [4], [5].

These datasets vary in the number of activities, sampling rate and types of sensors. This way, it is possible to examine the robustness of the methods on data with high variability.

The experimental results show that the currently used process to generate the data samples is not adequate to assess the quality of the activity recognition since current methods allow a data sample (or part of its content) appear in both training and testing (details in Section III-B), thereby, when appropriate data sample generation techniques are employed, the accuracy drops, on average, ten percentage points. Therefore, the results reported by previous works can be skewed and might not reflect their real performance.

To the best of our knowledge, this is the first study that tackles essential issues that compromise the results achieved in sensor-based human recognition, such issues have been responsible for skewing some previous results reported in the literature. We hope that this work allows a better comprehension of the challenges in human activity recognition based on wearable data and leads to further advancements of future works. All results of this work, including the data and implementations to reproduce them, are available¹.

II. HUMAN ACTIVITY RECOGNITION BASED ON WEARABLE SENSOR DATA

This section starts by describing some surveys related to the progress in human activity recognition based on wearable sensor data. Then, it introduces details regarding the works evaluated in our study, where we discuss approaches based on handcrafted and convolutional neural network.

A. Literature Reviews

One of the most comprehensive studies in human activity recognition based on wearable sensors is the work of Shoaib et al. [14]. Their work describes limitations and recommendations to online activity recognition using mobile phones. The term online refers to the implementation of the complete classification pipeline on the mobile phone, which consists on describing and classifying the signal. However, their work does not take into account convolutional neural network approaches, which nowadays are the most employed methods [5], [15]. On the other hand, Wang et al. [15] performed an extensive study regarding these approaches in the context of wearable sensors. In their work, the authors survey a number of deep learning based methods, including recurrent neural networks (RNN) and stacked autoencoders.

Stisen et al. [21] investigated the influence of heterogeneous devices on the final performance of the classifiers to perform activity recognition. For this purpose, the authors represented the activities using handcrafted features and employed popular classifiers such as nearest-neighbor, support vector machines and random forest. Additionally, in their work, the authors noticed that severe sampling instabilities occur in the devices, which contributes to a more challenging activity recognition. Similar to [14], [21], [15], Mukhopadhyay [22] performed

a detailed investigation regarding the advances in activity recognition associated with inertial data; however, focusing on the hardware context.

Different from these works, we do not summarize or review existing methods based on their reported results. Instead, we implement and conduct an extensive set of experiments on the methods to show crucial questions that affect them.

B. Methods based on Handcrafted Features

To represent the activities, Kwapisz et al. [7] extracted handcrafted features (i.e., average and standard deviation) from the raw signal. The authors also analyzed a set of classifiers to determinate the best one able to classify the categories of activities. With this purpose, Kwapisz et al. [7] examined multilayer perceptron, decision tree (J48) and logistic regression, where the first achieved the best classification results. Following these ideas, Catal et al. [8] proposed to apply ensemble techniques to combine these classifiers and compose the final predictor. By performing this process on the same features, Catal et al. [8] achieved a more accurate classification when compared with Kwapisz et al. [7].

To increase the discriminative power among activities, Kim et al. [9] divided an entire activity (in our context, a sample provided by temporal window process, details in Section IV-A) into a set of action units. Each action unit is represented by its average and correlation and is classified using bagging of decision trees. Finally, based on the proportion of each action unit, it is possible to predict which activity these action units belong. In contrast to [9], Kim et al. [10] proposed to use the boosting (compose of decision trees) with a smaller number of action units. The main difference between [9] and [10] is the number of action units and the classifier employed. Therefore, in this work, we only report the activity recognition accuracy of [9], which presented the best results in our experiments. Additionally, since the majority of the datasets do not provide enough information to build the activity units, we examine the work of [9] in terms of the features and classifier.

It is important to emphasize that some of the aforementioned works evaluate many classifiers. This way, to reduce the number of experiments and standardize their methods, our implementation considers only the best classifier of each method, according to their original paper.

C. Methods based on Convolutional Neural Networks

Another increasing line of research in human activity recognition based on wearable sensors aims at avoiding the design of handcrafted features, operation that requires human work and expert knowledge. These works employ convolutional neural networks (ConvNet) to learn the features and the classifier simultaneously.

Focusing on convolutional neural networks, Chen and Xue [3] employed a sophisticated ConvNet, where the input is taken from the raw signal (details in Section IV-A). They proposed a ConvNet architecture composed of three convolutional layers with 18, 36, 24 filters, respectively, followed by 2×1 max-pooling layers, each. To extract the association between two neighboring pairs of signal axes, at the first layer,

¹<http://www.sense.dcc.ufmg.br/activity-recognition-based-wearable-sensors/>

Chen and Xue applied a 12×2 (height \times width) convolutional kernel, while in the remaining layers the authors capture only the temporal relation with kernels 12×1 .

Similarly to [3], Jiang and Yin [11] introduced a ConvNet of two layers with convolutional kernels of 5×5 followed by 4×4 and 2×2 average-pooling layers, respectively. To improve the representation of the input data before presenting them to the ConvNet, the authors perform a process, referred to as *signal image*, which consists of following steps. Initially, a new signal is generated from a set of permutations using the axes of the raw signal. Then, a Fourier transform is applied to this new signal, producing the input to the ConvNet. Even with interesting results, the method proposed by Jiang and Yin [11] presents a remarkable drawback, the high computational cost since their method increases the input matrix size exponentially. This fact prevented us from applying their method in the dataset PAMAP2, as well as from conducting some of our experiments (more details in Section IV).

Following the hypothesis that different sensor modalities (i.e., accelerometer, gyroscope and magnetometer), should be convolved separately, Ha et al. [4] introduced a zero-padding between each heterogeneous modality to prevent them from being merged during the convolution process. Their architecture consists of two convolutional layers with 32 and 64 filters of 3×3 , respectively. However, due to this architecture, the heterogeneous modalities are convolved together at the second layer. To address this problem, Ha and Choi [5] suggested to introduce a zero-padding before starting the second layer of convolution. For this purpose, a zero-padding was inserted in the feature map generated by the first convolution layer so that the different modalities can be kept separated.

Ha and Choi [5] also demonstrated that ConvNets (2D convolutions) present better results compared with 1D convolutions (Conv1D), RNN or Long Short Term Memory networks (LSTM). In particular, even though recurrent-based networks have been successfully applied to speech recognition [23], [24] and natural language processing [25], [26], there exist few successful works that explore LSTMs in the context of human activity recognition based on wearable sensors [27], [28]. In general, recurrent-based networks have many hyperparameters to be set and do not present expressive results when compared to Conv2D. Thus, we do not contemplate this class of approaches in our experiments.

III. EVALUATION METHODOLOGY

The major concern of the research on human activity recognition based on wearable sensors is the lack of standard protocols to conduct experiments and report results. In other words, simple questions such as “*What is the evaluation metric to report the results?*”, “*How to generate the data samples from the raw signal?*” and “*What are the challenging datasets?*” have not been properly addressed in the existing works. The disregard of such questions prevents us from comparing the existing works and, as a consequence, it is not possible to determine the state-of-the-art in this task. For instance, while some works use *F-measure* to report the final performance of their methods, others employ accuracy. The problem becomes

worse when the authors choose different validation protocols. Moreover, an inadequate process to generate the data samples can bias the real performance of the methods (as it will be shown).

It is important to note that in other application areas, such as object detection [16], image classification [17] and face verification [18], [19], the aforementioned issues have been well-defined, which attracts more research efforts towards the application. Therefore, a standardized evaluation is an essential requirement for research in human activity recognition based on wearable sensor data.

Given this overview regarding the problems in wearable sensor data applied to activity recognition, the remaining of this section defines the evaluation methodology of this task. We start by describing the evaluation metrics and protocols employed by previous works to report the activity recognition performance in the context of wearable sensor data. Then, we discuss the traditional sample generation process and introduce the proposed approaches to perform this process.

A. Evaluation Metrics

There exists a set of metrics to measure the activity recognition performance, such as accuracy, recall, F-measure and Area Under the Curve (AUC). Table I summarizes the main evaluation metrics employed in human activity recognition based on wearable sensor data. Among the metrics listed in Table I, accuracy and F-measure are obvious choices. In particular, F-measure is more suitable since it is computed using the precision and recall, thereby, it is able to evaluate the activity recognition taking into account two different metrics.

B. Validation Protocols

An important step in recognition tasks is to separate the available data into training and testing sets. For this purpose, in the context of human activity recognition based on wearable sensors data, many works apply techniques such as *k*-fold cross validation, leave-one-subject-out, hold-out and leave-one-sample-out (a.k.a leave-one-out). The techniques of *k*-fold cross validation (with $k = 10$) and leave-one-subject-out are the traditional preferences, while few works employ the hold-out and leave-one-sample-out (this one due to the large number of executions), as showed in Table I.

We highlight that the leave-one-subject-out can be comprehended as a special case of the cross validation, where a subject can be seen as a fold, hence, the number of subjects determine the number of folds. Furthermore, it reflects a realistic scenario where a model is trained in an offline way [14] using the samples of some subjects and is tested with samples of unseen subjects. However, by using this protocol, the methods present high variance in accuracy from one subject to another since the same activity can be performed in different ways by the subjects.

C. Sample Generation Process

The first step to perform human activity recognition based on wearable sensor data is to generate the samples from the

TABLE I: Samples generation processes, metrics and validation protocol employed by different works. While it is clear that accuracy should be used as the metric, there is no sense in terms of validation protocol. *Unknown* denotes that the original paper does not report the technique employed.

Work	Data Generation	Evaluation Metrics	Validation Protocol
Pirttikangas [29]	Semi-Non-Overlapping-Window	Accuracy	4-fold cross validation
Suutala et al. [30]	Semi-Non-Overlapping-Window	Accuracy, Precision, Recall	4-fold cross validation
Kwapisz et al. [7]	Unknown	Accuracy	10-fold cross validation
Catal et al. [8]	Unknown	Accuracy, AUC, F-Measure	10-fold cross validation
Kim et al. [9]	Semi-Non-Overlapping-Window	F-measure	Unknown
Kim and Choi [10]	Semi-Non-Overlapping-Window	Accuracy, F-measure	Unknown
Chen and Xue [3]	Semi-Non-Overlapping-Window	Accuracy	Holdout
Jiang and Yin [11]	Unknown	Accuracy	Unknown
Ha et al. [4]	Semi-Non-Overlapping-Window	Accuracy	Hold out
Ha and Choi [5]	Semi-Non-Overlapping-Window	Accuracy	Leave-One-Subjet-Out
Yao et al. [27]	Semi-Non-Overlapping-Window	Accuracy.	Leave-One-Subjet-Out
Pan et al. [31]	Semi-Non-Overlapping-Window	Accuracy	Cross validation and Leave-One-Subjet-Out
Yang et al. [32]	Unknown	Accuracy, F-Measure	Hold out and Leave-One-Subjet-Out

raw signal. This process consists of splitting a raw signal into small windows of the same size, referred to as *temporal windows*. Then, the temporal windows from the signals are used as data samples, where they are split into training and test to learn and evaluate a model, respectively.

This section explains the process employed by previous works to generate the temporal windows, *Semi-Non-Overlapping-Window*. It also introduces two novel processes: *Full-Non-Overlapping-Window* and *Leave-One-Trial-Out*, both focuses on addressing the drawback of the existing process.

1) *Semi-Non-Overlapping-Window*: This is the most employed process to yield samples (temporal windows) for the activity recognition and it works as follows. Initially, the temporal sliding window technique (defined in Section IV-A) is applied on the raw signals, generating a set of data samples. Then, from these data samples, sets for training and test are created using some validation protocol (e.g., 10-fold cross validation). Since this process considers an overlap of 50% between windows, we called it of *Semi-Non-Overlapping-Window*.

A notable drawback of this process is that it is highly biased. This occurs because a window i and $i+1$ can appear in different folds of the cross validation (or any other protocol). Thereby, 50% of the content of these windows are equal because they present overlapping. As a consequence, 50% of a sample can appear in both training and testing at the same time, biasing the results. In other words, training and testing samples might be very close temporally, generating skewed results. Therefore, based on the second column of Table I, the results reported by the previous works do not reflect their real performance.

We emphasize that, in work work, the term *bias* refers to the fact that part of the sample’s content appears in the training and testing, simultaneously. According to the experiments, the methods drop the accuracy notably when changing from this process to another without this bias.

It is important to mention that on the leave-one-subject-out validation protocol, the semi-non-overlapping-window technique is not affected by bias, which is desirable, since the

samples of training and testing are separated by subjects. Therefore, the raw signal used to yield the samples (which can be temporally close) will either appear in the training or in the testing, but not in both.

2) *Full-Non-Overlapping-Window*: A simple way to handle the bias problem of the aforementioned process is to ensure that the windows have no temporal overlapping, guaranteeing that part of the window’s content does not appear in the training and testing, simultaneously. For this purpose, we propose the use of non-overlapping windows (overlap equal to zero between the temporal windows) to generate samples, process referred to as *Full-Non-Overlapping-Window*.

Even though the proposed full-non-overlapping-window process prevents the bias, it has the disadvantage of providing a reduced number of samples when compared to the semi-non-overlapping-window process (around 1.10 times fewer samples²) since the temporal windows no longer overlap.

3) *Leave-One-Trial-Out*: As we argued earlier, each process has a drawback that might cause a negative impact on the methods. For instance, semi-non-overlapping-window process produces biased results while the full-non-overlapping-window process generates few samples.

To avoid the aforementioned problems, we propose the *Leave-One-Trial-Out* process. A trial is the raw signal of one sequence of activities (or a single activity) performed by one subject. In this work, we propose to use the trials to ensure that samples generated by the same signal do not appear in the training and testing, simultaneously. To achieve this goal, we apply 10-fold cross validation on the raw signals, guaranteeing that the same trial (the same raw signal) does not appear in both training and testing. This is possible since k-fold cross validation ensures that the same sample (here, the trial - raw signal) appears in one fold only. Finally, for each fold, we generate the data samples (temporal windows) from raw signals, following the same process as in semi-non-overlapping-window. Throughout this process, we do not change the number of folds, this way, the number of folds is defined by k-fold cross validation (where k=10). It is important to mention that it is not possible/adequate to employ the trial

²Value computed using the average of all the datasets.

as a criterion to generate the folds since a single trial, in most cases, does not contain all the activities.

By using this process, we avoid: (1) bias since part of the window content never appear in the training and testing at the same time, and (2) small number of samples because the overlapping used in this process is the same employed in the semi-non-overlapping-window.

D. Standardization of the Wearable Sensors Datasets

Nowadays, there are many available datasets to perform human activity recognition based on wearable sensor data. These datasets present a wide range of sampling rate, number of activity categories and available sensors (see Table II), which enable us to evaluate the activity recognition in different scenarios. However, the lack of standardization of the captured data renders difficulties to develop a general framework able to perform activity recognition on all the datasets.

In general, the wearable sensors datasets can be divided into two groups with respect to the manner in which the activities were captured. The first group consists of activities where the user performs all the activities freely, i.e., there is no pause between the execution of one activity and the next one (e.g., MHEALTH, PAMAP2 and WISDM). The second group contains activities captured separately, i.e., a single activity is performed at a time (e.g., USC-HAD, UTD-MHAD and WHARF). This difference between these groups makes hard to perform a unified evaluation. Therefore, it is important to consolidated them into a single form. Intuitively, the first group of datasets can be converted into the second group, while the inverse is not possible, thereby, in this work we standardize all the datasets to simulate the second type of dataset.

As a final note, since the activity recognition datasets involve human participants, the ETHICS approval is required. This approval is of responsibility of the authors that proposed the datasets and can be found in the original works where the datasets were proposed.

IV. EXPERIMENTAL EVALUATION

We start this section by describing the experimental setup and the datasets in Section IV-A and IV-B, respectively. Afterwards, we present the experiments demonstrating the influence of the sample generation process on the activity recognition performance (Section IV-C). Then, we evaluate the impact of using subjects to separate the training and testing samples (Section IV-D), and investigate the activity recognition performance according to the datasets employed (Section IV-E). Finally, we compare the previous works using statistical tests in Section IV-F, and discuss their results in Section IV-G, where we define the state-of-the-art in activity recognition based on wearable data.

According to Table I, accuracy is the evaluation metric most employed by existing works, therefore, we have selected it to assess the activity recognition quality. To separate the data into training and testing (validation protocol), we used the 10-fold cross validation (since it is a common choice, as seen in Table I), except for the leave-one-subject-out protocol, where the folds are defined by the number of subjects.

A. Experimental Setup

1) *Temporal Sliding Window*: To increase the number of samples and enable the activity recognition to operate with a small latency (expected for real-time activity recognition), the works in the literature employ the temporal sliding window technique [10], [3], [4]. This technique consists of dividing the sample into subparts (windows) and considering each subpart as an entire activity. Specifically, each window becomes itself a sample that will be associated to a class label after its classification. A temporal sliding window can be defined as

$$w = [s_{k-t}, \dots, s_{k-2}, s_{k-1}, s_k]^\top, \quad (1)$$

where k represents the current signal captured by the sensor and t denotes the temporal sliding window size. The windows might overlap and the ones that do not fit within the temporal window are dropped. In other words, windows with the size smaller than t are discarded. Based on previous works [39], [40], we are using t equals to 5 seconds, which represents a good trade-off between the number of discarded samples and recognition accuracy.

2) *Convolutional Neural Network Setup*: Different from handcrafted approaches, where the authors provide enough information for reproducibility, most works based on ConvNets omit some important parameters, such as number of epochs, batch size and the optimizer used. To handle this problem and provide a fair comparison among this group of approaches, we set these parameters as follows. The maximum number of epochs was set as 200 and the method stops its training when the loss function reaches a value less or equal to 0.2. These values were set empirically by observing the trade-off between execution time and accuracy. Similarly, the batch size was set to 1000, except for the PAMAP2 dataset, where this value was of 250 due to memory issues. Finally, we employ the Adadelta optimizer [41] (except for the methods where the optimizer was specified by the author). In preliminary experiments, this optimizer presented the best accuracy when compared to SGD and RMSprop [42], besides providing an efficient execution time.

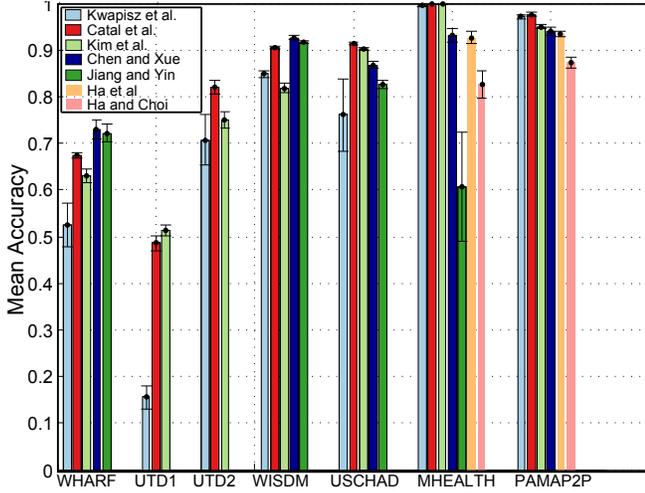
It is important to mention that the employment of deep architectures and large convolutional kernels makes impracticable the use of ConvNets in some datasets where the sampling rate is small, i.e., WHARF and WISDM (see Table II). In deep architectures, this occurs since the convolution process produces feature maps smaller than the input presented to it and its size can reach zero in deeper layers of a ConvNet. Therefore, it was not possible to execute some of the ConvNets considered in this work for all datasets.

B. Datasets

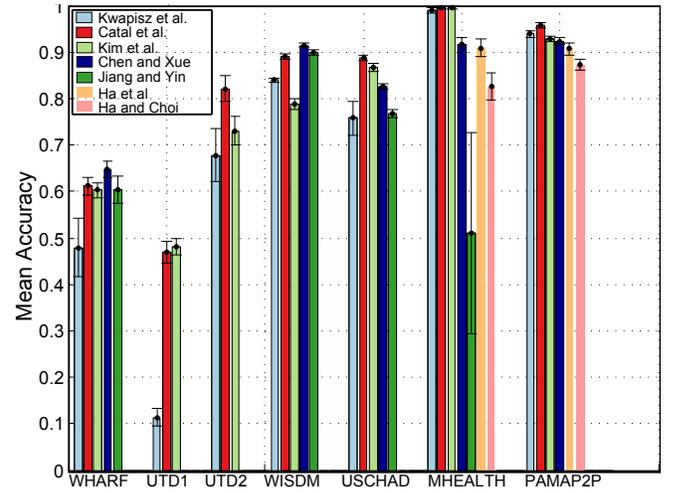
The datasets evaluated in our study consider a variety of sampling rate, number of activities categories and degree of difficulty. To select these datasets, we consider the ones which provide enough information regarding the capturing of the data, activities, subjects and the employed sensors. Additionally, we label a dataset as unbalanced when its largest class has four times more samples than the smallest class.

TABLE II: Main features of the datasets used in this work. The number of samples was computed using the semi-non-overlapping-window process with leave-one-subject-out validation protocol. Acc, Gyro, Mag and Temp indicate accelerometer, gyroscope, magnetometer and temperature, in this order.

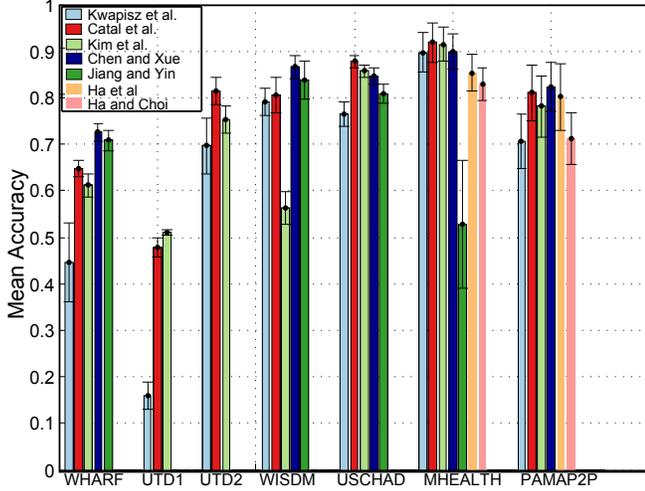
Dataset	Frequency (Hz)	#Sensors	#Activities	#Subjects	#Trials	#Samples	Balanced
MHEALTH [33]	50	3 (Acc, Gyro, Mag)	12	10	262	2555	True
PAMAP2 [34]	100	4 (Acc, Gyro, Mag, Temp)	12	10	108	7522	False
USCHAD [35]	100	2 (Acc, Gyro)	12	15	840	9824	False
UTD-1 [36]	50	2 (Acc, Gyro)	21	9	617	3771	True
UTD-2 [36]	50	2 (Acc, Gyro)	5	9	190	1137	True
WHARF [37]	32	1 (Acc)	12	17	884	3871	False
WISDM [38]	20	1 (Acc)	7	36	402	20846	False



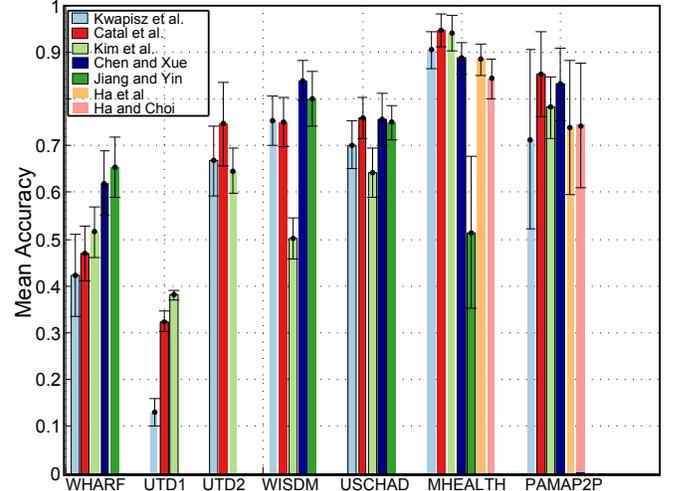
(a) Results using the semi-non-overlapping-window and cross validation (SNVC) combination.



(b) Results using the full-non-overlapping-window and cross validation (FNVC) combination.



(c) Results using the leave-one-trial-out and cross validation (LTCV) combination



(d) Results using the semi-non-overlapping-window and leave-one-subject-out (SNLS) combination

Fig. 1: Mean accuracy achieved by the methods using different process to generate the data samples (best visualized in color). Black bars denote the confidence interval.

The main features of the datasets are summarized in Table II and briefly described as follows. The MHEALTH and PAMAP2 datasets consist of activities captured from sensors placed on the subject's chest, wrist and ankle. Similarly, the activities of the WHARF dataset were obtained from a sensor placed on the subject's wrist. Different from these datasets, the activities of the WISDM dataset were captured using a

sensor located in the user's pocket focusing on convenience and comfort for the subject. Finally, the activities of the UTD-MHAD dataset were captured using two configurations, the sensor placed at the subject's wrist and subject's thigh, UTD-1 and UTD-2, in this order.

C. Sample Generation Processes and Validation Protocols

This experiment intends to demonstrate that there is a considerable variance in the results achieved by the methods when different sample generation processes and validation protocols are considered. In particular, we can use different combinations of sample generation processes and validation protocols to perform activity recognition, as seen in Table III. However, to reduce the number of experiments, we use the proposed sample generation processes on the 10-fold cross validation only, which is the most employed validation protocol, reducing the number of possible combinations to four. In addition, once selected a combination of Table III, we ensure that the training and testing samples are the same to all the methods. In this way, we provide an adequate and fair comparison.

Figure 1 shows the mean accuracy and the confidence interval of the methods when evaluated on different combinations of sample generation processes and validation protocols (see Table III). According to Figure 1(a), it is possible to note that the combination SNCV reports the highest mean accuracy when compared to the other combinations. This happens due to bias produced by its sample generation process, where the content of a window can appear in both training and testing. Therefore, the works that employ the SNCV combination might have their results highly skewed.

Figure 1(b) shows the results of the methods using the FNCV combination. According to the results, the SNCV and FNCV achieve similar mean accuracy, where the second one presents, slightly, inferior results. This effect is expected since FNCV also has a sampling bias because the trial used to generate the samples (which do not present overlapping) can produce samples of training and testing at the same time. In other words, parts of the same trial can appear in training and testing. On the other hand, when the LTCV combination is considered (Figure 1(c)), the mean accuracy drops significantly. This is a consequence of the data generation, which do not have any type of bias. It is possible to note this behavior by observing the results achieved in the PAMAP2 and MHEALTH datasets, where the methods had its accuracy reduced drastically when compared to the results of the SNCV and FNCV combinations.

Figure 1(d) illustrates the results achieved on the SNLS combination, where the mean accuracy had the smallest performance. This occurs since SNLS is invariant to bias, since the samples of training and testing are separated by subjects,

TABLE III: Combinations between samples generation processes and validation protocols. The symbol '-' denotes we do not consider the respective combination.

	10-Fold Cross Validation	Leave-One Subject-Out
Semi-Non Overlapping-Window	SNCV	SNLS
Full-Non Overlapping-Window	FNCV	-
Leave-One Trial-Out	LTCV	-

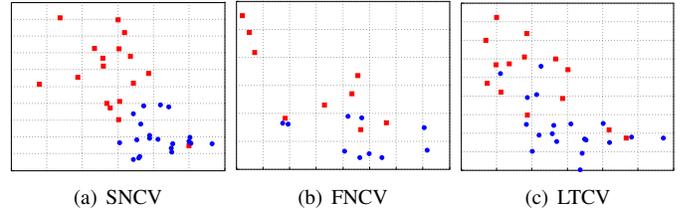


Fig. 2: Feature space provided by the two first components of Linear Discriminant Analysis on the features suggested by [7]. Note that from Figure (a) to (c) the separability of the classes is decreasing, which is a consequence of the process used to generate the data samples. To provide a better visualization, we show only two activities categories (best visualized in color).

which means the raw signal used to yield the samples either will appear in the training or testing only, causing an effect similar to the proposed leave-one-trial-out.

Based on the aforementioned discussion, it is possible to note that there exists a high variance in the results depending on the combinations between samples generation processes and validation protocols. To demonstrate this, let us compare the number of methods which achieved a mean accuracy above 80%³ to a determinate data sample generation process, by considering all the datasets. From SNCV to FNCV (Figure 1(a) and (b)) and FNCV to LTCV (Figure 1(b) and (c)), this number decreased from 21 to 20 and from 20 to 16, in this order. These values indicate that the activity recognition becomes more difficult with respect to the samples generation process employed. In particular, this variance in the results is a direct effect of the bias introduced during the process of generating the data samples, where part of the window's content appear both in training and testing. This remark is easier to notice when we compare SNCV and SNLS, Figures 1(a) and (d). From this comparison, the number of methods which achieved an accuracy above 80% decreased from 21 to 10, when using SNLS instead of SNCV.

To illustrate the data behavior when the samples generation process is changed, we set the validation protocol (10-fold cross validation) and vary the samples generation process, then, we project the training samples onto the two first components of Linear Discriminant Analysis (LDA) [43], Figure 2. From this figure, it is clear that the class separability decreases according to the process employed.

Another point to be evaluated regarding the combinations in Table III is the accuracy variance between the folds of the validation protocol. To that, we select the method of Catal et al. [8] and report their accuracy obtained for each fold using the PAMAP2 dataset⁴, Figure 3. Figure 3 shows that SNCV and FNCV present minimal variance. However, it increases when the methods are evaluated on the LTCV and SNLS combinations.

Based on the experiments conducted in this section, we

³We select this value empirically just for discussion.

⁴Specifically, we choose this method since it presents the best performance and we select the PAMAP2 dataset because it provides exactly 10 subjects, which enables us to compare the i -th subject with the i -th fold.

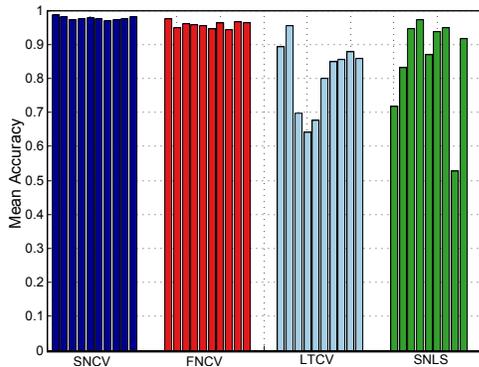


Fig. 3: Accuracy obtained by each fold using different combinations of validation protocols and samples generation processes (best visualized in color).

demonstrated that the methods have their performance extremely associated with the combination between samples generation process and validation protocol. For instance, semi-non-overlapping-window using cross validation introduces bias, which is inadequate for conducting experiments, while semi-non-overlapping-window using leave-one-subject-out is bias-invariant. Unfortunately, according to Table I, the majority of the methods employ the first combination. On the other hand, we showed that this bias can be slightly reduced or completely removed with the employment of our proposed sample generation process.

D. Limitations of Leave-One-Subject-Out

In this experiment, we intend to show the drawback of leave-one-subject-out validation protocol and propose an alternative. According to Figure 1(d) and Table IX, the confidence interval (which is directly affected by the accuracy variance of the folds) is extremely high, as indicated in datasets such as WHARF, MHEALTH and PAMPA2P. This behavior is due to high variability in the samples provided by the folds (recall that for leave-one-subject-out, folds are the subjects). Specifically, different subjects can perform the activities in distinct ways; hence, the learned model is not able to produce enough generalization to correctly classifying the testing subject. A drawback regarding this issue is that the methods become statistically equivalent (since the confidence intervals are large and overlap each other [20]).

To face this problem, we propose SNLS \times 10. The idea behind SNLS \times 10 is to measure the accuracy and the confidence interval by using the variability from the training samples, which is small, instead of the variability from the subjects, which is high causing a large confidence interval. With this purpose, we execute the leave-one-subject-out 10 times and for each turn we select just 80% of the training samples (the same samples used in traditional SNLS) to learn the model.

Figure 4 shows that performing SNLS \times 10 decreases the confidence interval when compared to the traditional SNLS, enabling us to compare the methods using the confidence interval. According to Figure 4, methods of Catal et al. [8] and Kim et al. [9] achieved the Top 1 best accuracy besides being

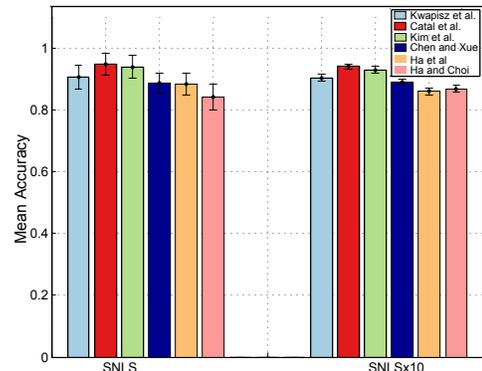


Fig. 4: Comparison of the confidence interval between SNLS and SNLS \times 10 (best visualized in color).

statistically different from the remaining methods. The Top 2 best accuracy was obtained by the methods of Kwapisz et al. [7] and Chen and Xue [3], which are statistically equivalent. Finally, the Top 3 best result was achieved by Ha et al. [4] and Ha and Choi [5].

Note that it is not possible to list the rank above using the traditional SNLS combination since all the methods (except [5]) are statistically equivalent. However, the computational cost of SNLS \times 10 increases because it is necessary to execute more turns of the leave-one-subject-out (\times 10). This fact forced us to perform the SNLS \times 10 only in the MHEALTH dataset where it is possible to execute all the methods.⁵

According to the results achieved in this experiment, we conclude that, though SNLS is an adequate combination to avoid bias and to force the model to be generalized, it is unsuitable to perform statistical tests among the methods.

E. Comparison of Datasets

This experiment intends to determine which are the most challenging datasets in human activity recognition based on wearable sensor data. As we argued in previous experiments, semi-non-overlapping-window and full-non-overlapping-window when combined with cross validation, are not suitable to conduct experiments. Therefore, in this experiment, our discussion is conducted considering the combinations of leave-one-trial-out with cross validation and semi-non-overlapping-window with leave-one-subject-out (LTCV and SNLS, in this order, as defined in Table III).

In this experiment, we measure the mean accuracy obtained by all the methods on each dataset (last row in Tables VIII and IX). According to the results, the most challenging dataset is UTD-1, with a mean accuracy of 38.25% and 27.84% when considering the LTCV and SNLS combinations, respectively. We believe that the low accuracy is due to two main reasons: the low sampling rate and the large number of activities. In particular, the large number of activities to be recognized in UTD-1 makes more challenging the recognition as can be observed by the contrast in accuracy regarding UTD-2, where

⁵Due to the high number of executions and memory constraints, in this experiment was not possible to execute the method of Jiang and Yin [11].

TABLE IV: Number of times where a method was statistically equivalent to other.

	MHEALTH		PAMAP2	
	LTCV	SNLS	LTCV	SNLS
Kwapisz et al. [7]	4	4	2	5
Catal et al. [8]	3	2	3	5
Kim et al. [9]	3	2	5	5
Chen and Xue [3]	4	3	3	5
Jiang and Yin [11]	0	0	0	0
Ha et al. [4]	3	3	3	5
Ha and Choi [5]	1	1	2	5

there are a smaller number of activities. Similarly, on the WHARF dataset, the methods also achieved a low performance, with a mean accuracy of 62.76% and 53.55% when employing the LTCV and SNLS combinations, respectively. On the other hand, the methods achieved the best performance on the MHEALTH dataset, where the mean accuracy was superior to 83%. This is an effect of the high sampling rate and the number of available sensors provided by this dataset, which make easier the recognition of the activities.

Observe that our discussion above considers the mean accuracy of all methods, which can be skewed by single methods, for instance, the methods of Kim et al. [9] and Jiang and Yin [11] on the WISDM and MHEALTH datasets, respectively. However, by examining the datasets using only the best method, these claims still remain valid.

F. Statistical Evaluation

This experiment intends to show whether the performance achieved by works in human activity recognition is, de-facto, statistically superior or equivalent to other. In addition, it focuses on showing that the proposed leave-one-trial-out with cross validation (LTCV) is more suitable than the combination semi-non-overlapping-window with leave-one-subject-out (SNLS) for conducting statistical tests.

Following Jain [20], we perform a *unpaired t-test*, which works as follows. For each pair of methods, we computed the confidence interval, using a confidence of 90%, from the difference between their mean accuracies. In the cases where the resulting confidence interval includes the zero value, the methods are statistically equivalent. Otherwise, they are statistically different. Note that this process is the same suggested in [20]. Due to the high number of comparisons, we perform this evaluation on MHEALTH and PAMAP2 datasets only, which are the ones where it is possible to execute a larger number of works. Table IV shows the number of times where a method was statistically equivalent to another using the unpaired t-test.

For both datasets and combinations evaluated, the method of Jiang and Yin [11] is the one with the lower number of statistical equivalence. However, sometimes, their accuracy was statistically inferior. On the other hand, on the PAMAP2 dataset and the SNLS combination, the methods of Ha et al. [4] and, Ha and Choi [5] were the ones with the smaller number of statistical equivalence.

According to Table IV, a considerable amount of methods shown to be statistically equivalents. This is an effect of the

large confidence interval (denoted by brackets in Tables VI-IX and black bars in Figure 1), caused by the high variance in accuracy, as seen in Figure 3. In particular, the LTCV combination has a smaller variance than SNLS (0.0107 against 0.1448, respectively). As a consequence, LTCV presented a smaller number of methods which are statistically equivalents (mainly on the PAMAP2 dataset). This evidence shows that our leave-one-trial-out when using cross validation is adequate to conduct statistical tests.

G. The State-of-The-Art

Our last experiment focuses on defining the state-of-the-art methods in human activity recognition based on wearable sensor data. To this end, following previous experiments, we discuss the results using the combination of LTCV and SNLS. In addition, since most of the methods are statistically equivalent, as shown in the earlier experiment, we are not considering statistically difference to determinate if a method is superior than another.

Based on the results presented in Tables VIII and IX, we report the number of datasets where the methods achieved the best (Top 1) accuracy. Additionally, since there exists a high variance in accuracy among datasets, as seen in Tables VIII-IX and Figure 1, we consider the second (Top 2) and the third (Top 3) best accuracy, aiming to make a fairer comparison. Table V summarizes these results.

Regarding the methods based on ConvNets, the approach with the best performance is the method of Chen and Xue [3], achieving the Top1 one and three times when evaluated on LTCV and SNLS, respectively. This result is a consequence of the filter shapes, which are adequate to capture the temporal and spatial pattern of the signal. On the other hand, the methods of Jiang and Yin [11], Ha et al. [4] and, Ha and Choi [5] were not able to achieve good results. We believe that these inaccurate results are an effect of their convolutional filters (3×3 and 5×5), which capture a small temporal pattern besides being sensitive to noise. In particular, the methods [4] and [5] were evaluated only on two datasets (due to issues of the network architecture). However, by normalizing their results by the number of datasets evaluated, they still do not present good performance. Finally, by considering handcrafted and ConvNets approaches, the more accurate methods are the approaches of Catal et al. [8] and Chen and Xue [3]. This result indicates that, though ConvNets-based approaches have

TABLE V: Number of datasets where a method achieved the first (Top1), second (Top2) and third (Top3) best performance. Values after '/' denote the number of datasets where were possible to execute the method.

Method	LTCV			SNLS		
	Top1	Top2	Top3	Top1	Top2	Top3
Kwapisz et al. [7]	0/7	0/7	1/7	0/7	2/7	0/7
Catal et al. [8]	3/7	2/7	2/7	4/7	1/7	0/7
Kim et al. [9]	1/7	3/7	1/7	1/7	1/7	3/7
Chen and Xue [3]	3/5	0/5	2/5	1/5	3/5	0/5
Jiang and Yin [11]	0/4	2/4	0/4	1/4	1/4	1/4
Ha et al. [4]	0/2	0/2	1/2	0/2	0/2	0/2
Ha and Choi [5]	0/2	0/2	0/2	0/2	0/2	0/2

TABLE VI: Mean accuracy and confidence interval using the Semi-Non-Overlapping-Window and Cross Validation (SNCV) combination. The symbol '-' denotes which was not possible to execute the method on the respective dataset.

	MHEALTH	PAMAP2	USCHAD	UTD-1	UTD-2	WHARF	WISDM	Mean Accuracy
Kwapisz et al. [7]	99.49 [99.28, 99.71]	97.09 [96.72, 97.46]	76.08 [68.27, 83.89]	15.54 [13.11, 17.96]	70.73 [65.32, 76.13]	52.47 [47.70, 57.25]	84.83 [84.08, 85.58]	70.89
Catal et al. [8]	99.92 [99.83, 1.000]	97.66 [97.34, 97.98]	91.37 [91.00, 91.73]	48.57 [47.01, 50.12]	82.00 [80.55, 83.45]	67.35 [66.79, 67.90]	90.62 [90.32, 90.92]	82.49
Kim et al. [9]	99.88 [99.77, 99.99]	94.99 [94.54, 95.44]	90.21 [89.84, 90.58]	51.42 [50.29, 52.55]	74.97 [73.12, 76.82]	62.91 [61.41, 64.41]	81.75 [80.71, 82.79]	79.44
Chen and Xue [3]	93.19 [91.75, 94.62]	94.03 [93.12, 94.95]	86.78 [86.00, 87.55]	—	—	72.92 [70.80, 75.03]	92.65 [92.18, 93.12]	87.91
Jiang and Yin [11]	60.69 [48.88, 72.50]	—	82.55 [81.73, 83.36]	—	—	72.22 [70.47, 73.98]	91.68 [91.43, 91.94]	76.78
Ha et al. [4]	92.60 [91.29, 93.91]	93.42 [92.74, 94.10]	—	—	—	—	—	93.00
Ha and Choi [5]	84.46 [82.42, 86.51]	90.68 [89.93, 91.43]	—	—	—	—	—	92.33
Mean Accuracy	90.03	94.64	85.39	38.51	75.90	65.57	88.30	

TABLE VII: Mean accuracy and confidence interval using the Full-Non-Overlapping-Window and Cross Validation (FNVCV) combination. The symbol '-' denotes which was not possible to execute the method on the respective dataset.

	MHEALTH	PAMAP2	USCHAD	UTD-1	UTD-2	WHARF	WISDM	Mean Accuracy
Kwapisz et al. [7]	99.03 [98.41, 99.64]	93.86 [93.03, 94.70]	75.74 [72.06, 79.43]	11.34 [09.42, 13.27]	67.83 [62.21, 73.45]	47.92 [41.52, 54.31]	83.93 [83.49, 84.37]	68.52
Catal et al. [8]	99.63 [99.32, 99.93]	95.77 [95.23, 96.31]	88.76 [88.12, 89.41]	46.90 [44.56, 49.23]	82.10 [79.28, 84.93]	61.15 [59.35, 62.96]	89.03 [88.48, 89.57]	80.47
Kim et al. [9]	99.70 [99.40, 100.0]	92.84 [92.31, 93.37]	86.80 [85.89, 87.71]	48.02 [46.28, 49.76]	73.03 [69.91, 76.16]	60.31 [58.78, 61.85]	78.82 [77.72, 79.93]	77.07
Chen and Xue [3]	91.61 [89.98, 93.25]	92.36 [91.55, 93.18]	82.63 [82.04, 83.23]	—	—	64.77 [62.95, 66.59]	91.47 [91.06, 91.88]	84.56
Jiang and Yin [11]	51.00 [29.29, 72.70]	—	76.77 [75.78, 77.77]	—	—	60.28 [57.37, 63.20]	90.02 [89.47, 90.58]	69.51
Ha et al. [4]	90.85 [88.96, 92.75]	90.77 [89.45, 92.09]	—	—	—	—	—	90.81
Ha and Choi [5]	82.70 [79.76, 85.64]	87.36 [86.20, 88.52]	—	—	—	—	—	85.03
Mean Accuracy	87.78	92.16	82.14	35.41	74.32	58.85	86.65	

TABLE VIII: Mean accuracy and confidence interval using the Leave-One-Trial-Out and Cross Validation (LTCV) combination. The symbol '-' denotes which was not possible to execute the method on the respective dataset.

	MHEALTH	PAMAP2	USCHAD	UTD-1	UTD-2	WHARF	WISDM	Mean Accuracy
Kwapisz et al. [7]	89.75 [85.52, 93.98]	70.58 [64.74, 76.41]	76.52 [73.97, 79.07]	15.99 [13.00, 18.97]	69.61 [63.68, 75.54]	44.51 [36.00, 53.02]	79.08 [76.26, 81.90]	63.71
Catal et al. [8]	91.84 [87.67, 96.01]	81.03 [75.02, 87.04]	87.77 [86.52, 89.02]	47.80 [45.70, 49.89]	81.37 [78.43, 84.32]	64.84 [63.05, 66.63]	80.52 [76.66, 84.38]	76.45
Kim et al. [9]	91.51 [87.95, 95.06]	78.08 [71.63, 84.54]	85.70 [84.28, 87.12]	50.98 [50.45, 51.51]	75.27 [72.23, 78.31]	61.12 [58.55, 63.68]	56.26 [52.76, 59.76]	71.27
Chen and Xue [3]	89.95 [86.21, 93.70]	82.32 [77.14, 87.50]	84.66 [83.04, 86.29]	—	—	72.55 [70.57, 74.54]	86.55 [84.14, 88.96]	83.20
Jiang and Yin [11]	52.78 [39.05, 66.52]	—	80.73 [78.70, 82.75]	—	—	70.79 [68.69, 72.88]	83.82 [79.68, 87.96]	72.03
Ha et al. [4]	85.31 [81.43, 89.20]	80.13 [72.99, 87.27]	—	—	—	—	—	82.71
Ha and Choi [5]	82.75 [79.23, 86.26]	71.19 [65.70, 76.69]	—	—	—	—	—	76.96
Mean Accuracy	83.41	77.22	83.07	38.25	75.41	62.76	77.24	

TABLE IX: Mean accuracy and confidence interval using the Semi-Non-Overlapping-Window and Leave-One-Subject-Out (SNLS) combination. The symbol '-' denotes which was not possible to execute the method on the respective dataset.

	MHEALTH	PAMAP2	USCHAD	UTD-1	UTD-2	WHARF	WISDM	Mean Accuracy
Kwapisz et al. [7]	90.41 [86.54, 94.28]	71.27 [52.07, 90.47]	70.15 [65.06, 75.24]	13.04 [10.19, 15.90]	66.67 [59.21, 74.14]	42.19 [33.40, 50.98]	75.31 [70.07, 80.55]	61.29
Catal et al. [8]	94.66 [91.17, 98.15]	85.25 [76.27, 94.22]	75.89 [71.62, 80.16]	32.45 [30.18, 34.71]	74.67 [65.75, 83.58]	46.84 [41.02, 52.67]	74.96 [69.66, 80.27]	69.29
Kim et al. [9]	93.90 [90.05, 97.74]	81.57 [73.50, 89.64]	64.20 [58.88, 69.53]	38.05 [37.01, 39.09]	64.60 [59.73, 69.47]	51.48 [46.11, 56.84]	50.22 [45.85, 54.59]	63.43
Chen and Xue [3]	88.67 [85.38, 91.96]	83.06 [75.40, 90.71]	75.58 [70.05, 81.11]	—	—	61.94 [55.02, 68.86]	83.89 [79.72, 88.06]	78.62
Jiang and Yin [11]	51.46 [35.35, 67.57]	—	74.88 [71.28, 78.48]	—	—	65.35 [58.81, 71.88]	79.97 [74.21, 85.73]	67.91
Ha et al. [4]	88.34 [84.91, 91.78]	73.79 [59.56, 88.02]	—	—	—	—	—	81.06
Ha and Choi [5]	84.23 [80.01, 88.44]	74.21 [60.91, 87.50]	—	—	—	—	—	79.21
Mean Accuracy	84.52	78.19	72.14	27.84	68.64	53.55	72.87	

presented remarkable results in human activity recognition based on wearable data, handcrafted approaches are able to achieve comparable results.

V. CONCLUSIONS

This work conducted an extensive set of experiments to demonstrate essential issues which currently are not considered during the evaluation of the human activity recognition based on wearable sensor data. The main issue is regarding the process employed to generate the data samples, where the traditional process is susceptible to bias leading to skewed results. To demonstrate this, we investigate novel techniques to generate the data samples, which focus on reducing and removing this bias. According to our experiments, the accuracy drops considerably when appropriated data generation processes (bias-invariant) are used. Hence, the results reported by previous works can be skewed and do not reflect their real performance. In addition, throughout the experiments, we implement several top-performance methods and evaluated them on many popular and publicly available datasets. Thereby, we define the state-of-the-art methods in human activity recognition based on wearable sensor data.

We highlight that, different from previous studies and surveys, our work does not summarize or discuss existing

methods based on their reported results, which makes this work, to the best of our knowledge, the first that implements, groups and handles important issues regarding the activity recognition associated with wearable sensor data.

ACKNOWLEDGMENTS

The authors would like to thank the Brazilian National Research Council - CNPq (Grant #311053/2016-5), the Minas Gerais Research Foundation - FAPEMIG (Grants APQ-00567-14 and PPM-00540-17) and the Coordination for the Improvement of Higher Education Personnel - CAPES (DeepEyes Project). Part of the results presented in this paper were obtained through research on a project titled "HAR-HEALTH: Reconhecimento de Atividades Humanas associadas a Doenças Crônicas", sponsored by Samsung Eletrônica da Amazônia Ltda. under the terms of Brazilian federal law No. 8.248/91.

REFERENCES

- [1] A. M. Khan, Y. K. Lee, S. Y. Lee, and T. S. Kim, "A triaxial accelerometer-based physical-activity recognition via augmented-signal features and a hierarchical recognizer," *IEEE Transactions on Information Technology in Biomedicine*, vol. 14, pp. 1166–1172, 2010.
- [2] Akram Bayat, Marc Pomplun, and Duc A. Tran, "A study on human activity recognition using accelerometer data from smartphones," in *Conference on Future Networks and Communications*, 2014, pp. 450–457.

- [3] Yuqing Chen and Yang Xue, "A Deep Learning Approach to Human Activity Recognition Based on Single Accelerometer," in *IEEE International Conference on Systems, Man, and Cybernetics*, 2015, pp. 1488–1492.
- [4] Sojeong Ha, Jeong-Min Yun, and Seungjin Choi, "Multi-modal convolutional neural networks for activity recognition," in *IEEE International Conference on Systems, Man, and Cybernetics*, 2015, pp. 3017–3022.
- [5] Sojeong Ha and Seungjin Choi, "Convolutional neural networks for human activity recognition using multiple accelerometer and gyroscope sensors," in *International Joint Conference on Neural Networks*, 2016, pp. 381–388.
- [6] Ming Zeng, Le T. Nguyen, Bo Yu, Ole J. Mengshoel, Jiang Zhu, Pang Wu, and Joy Zhang, "Convolutional neural networks for human activity recognition using mobile sensors," in *6th International Conference on Mobile Computing, Applications and Services*, 2014, pp. 197–205.
- [7] Jennifer R. Kwapisz, Gary M. Weiss, and Samuel Moore, "Activity recognition using cell phone accelerometers," *SIGKDD Explorations*, vol. 12, no. 2, pp. 74–82, 2010.
- [8] Cagatay Catal, Selin Tufekci, Elif Pirmitt, and Guner Kocabag, "On the use of ensemble of classifiers for accelerometer-based activity recognition," *Applied Soft Computing*, vol. 37, pp. 1018–1022, 2015.
- [9] Hyun-Jun Kim, Mira Kim, Sunjae Lee, and Young Sang Choi, "An Analysis of Eating Activities for Automatic Food Type Recognition," in *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*, 2012, pp. 1–5.
- [10] Hyun-Jun Kim and Young Sang Choi, "Eating Activity Recognition for Health and Wellness: A Case Study on Asian Eating Style," in *IEEE International Conference on Consumer Electronics*, 2013, pp. 446–447.
- [11] Wenchao Jiang and Zhaozheng Yin, "Human activity recognition using wearable sensors by deep convolutional neural networks," in *23rd Annual Conference on Multimedia Conference*, 2015, pp. 1307–1310.
- [12] Artur Jordao, Leonardo Antônio Borges Torres, and William Robson Schwartz, "Novel approaches to human activity recognition based on accelerometer data," *Signal, Image and Video Processing*, pp. 1–8, 2018.
- [13] Oscar D. Lara and Miguel A. Labrador, "A survey on human activity recognition using wearable sensors," *IEEE Communications Surveys and Tutorials*, vol. 15, pp. 1192–1209, 2013.
- [14] Muhammad Shoaib, Stephan Bosch, Ozlem Incel, Hans Scholten, and Paul Havinga, "A Survey of Online Activity Recognition using Mobile Phones," *Sensors*, vol. 15, pp. 2059–2085, 2015.
- [15] Jindong Wang, Yiqiang Chen, Shuji Hao, Xiaohui Peng, and Lisha Hu, "Deep learning for sensor-based activity recognition: A survey," *CoRR*, vol. abs/1707.03502, 2017.
- [16] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results," <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [17] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, vol. 115, pp. 211–252, 2015.
- [18] Gary B. Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," Tech. Rep., University of Massachusetts, Amherst, 2007.
- [19] Lior Wolf, Tal Hassner, and Itay Maoz, "Face recognition in unconstrained videos with matched background similarity," in *Conference on Computer Vision and Pattern Recognition*, 2011, pp. 529–534.
- [20] Raj Jain, *The art of computer systems performance analysis: techniques for experimental design, measurement, simulation, and modeling*, Wiley professional computing. John Wiley & Sons, 1990.
- [21] Allan Stisen, Henrik Blunck, Sourav Bhattacharya, Thor Siiger Prentow, Mikkel Baun Kjærgaard, Anind K. Dey, Tobias Sonne, and Mads Møller Jensen, "Smart devices are different: Assessing and mitigating mobile sensing heterogeneities for activity recognition," in *Conference on Embedded Networked Sensor Systems*, 2015, pp. 127–140.
- [22] S. C. Mukhopadhyay, "Wearable sensors for human activity monitoring: A review," *IEEE Sensors Journal*, vol. 15, pp. 1321–1330, 2015.
- [23] Alex Graves, Abdel-rahman Mohamed, and Geoffrey E. Hinton, "Speech recognition with deep recurrent neural networks," in *International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 6645–6649.
- [24] Klaus Greff, Rupesh Kumar Srivastava, Jan Koutnik, Bas R. Steunebrink, and Jürgen Schmidhuber, "LSTM: A search space odyssey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 10, pp. 2222–2232.
- [25] Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Pei-hao Su, David Vandyke, and Steve J. Young, "Semantically conditioned lstm-based natural language generation for spoken dialogue systems," in *Conference on Empirical Methods in Natural Language Processing*, 2015, pp. 1711–1721.
- [26] Shuohang Wang and Jing Jiang, "Learning natural language inference with LSTM," in *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016, pp. 1442–1451.
- [27] Shuochao Yao, Shaohan Hu, Yiran Zhao, Aston Zhang, and Tarek F. Abdelzaher, "Deepsense: A unified deep learning framework for time-series mobile sensing data processing," pp. 351–360, 2017.
- [28] Abdulmajid Murad and Jae-Young Pyun, "Deep recurrent neural networks for human activity recognition," *Sensors*, vol. 17, pp. 2556, 2017.
- [29] Susanna Pirttikangas, Kaori Fujinami, and Tatsuo Nakajima, "Feature selection and activity recognition from wearable sensors," in *Ubiquitous Computing Systems*, 2006, pp. 516–527.
- [30] Jaakko Suutala, Susanna Pirttikangas, and Juha Röning, "Discriminative temporal smoothing for activity recognition from wearable sensors," in *Ubiquitous Computing Systems*, 2007, pp. 182–195.
- [31] Madhuri Panwar, S. Ram Dyuthi, K. Chandra Prakash, Dwaipayan Biswas, Amit Acharyya, Koushik Maharatna, Arvind Gautam, and Ganesh R. Naik, "CNN based approach for activity recognition using a wrist-worn accelerometer," in *Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, 2017, pp. 2438–2441.
- [32] Jianbo Yang, Minh Nhut Nguyen, Phyo Phyo San, Xiaoli Li, and Shonali Krishnaswamy, "Deep convolutional neural networks on multichannel time series for human activity recognition," in *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015, pp. 3995–4001.
- [33] Oresti Baños, Rafael García, Juan A. Holgado-Terriza, Miguel Damas, Héctor Pomares, Ignacio Rojas Ruiz, Alejandro Saez, and Claudia Villalonga, "mhealthroid: A novel framework for agile development of mobile health applications," in *Ambient Assisted Living and Daily Activities - 6th International Work-Conference*, 2014, pp. 91–98.
- [34] Attila Reiss and Didier Stricker, "Introducing a new benchmarked dataset for activity monitoring," in *16th International Symposium on Wearable Computers*, 2012, pp. 108–109.
- [35] Mi Zhang and Alexander A. Sawchuk, "Usc-had: A daily activity dataset for ubiquitous activity recognition using wearable sensors," in *ACM Conference on Ubiquitous Computing*, 2012, pp. 1036–1043.
- [36] Chen Chen, Roozbeh Jafari, and Nasser Kehtarnavaz, "UTD-MHAD: A multimodal dataset for human action recognition utilizing a depth camera and a wearable inertial sensor," in *IEEE International Conference on Image Processing*, 2015, pp. 168–172.
- [37] Barbara Bruno, Fulvio Mastrogiovanni, and Antonio Sgorbissa, "Wearable Inertial Sensors: Applications, Challenges, and Public Test Benches," *IEEE Robotics and Automation Magazine*, vol. 22, pp. 116–124, 2015.
- [38] Jeffrey W. Lockhart, Gary M. Weiss, Jack C. Xue, Shaun T. Gallagher, Andrew B. Grosner, and Tony T. Pulickal, "Design considerations for the wisdom smart phone-based sensor mining architecture," in *Fifth International Workshop on Knowledge Discovery from Sensor Data*, 2011, pp. 25–33.
- [39] Dan Morris, T. Scott Saponas, Andrew Guillory, and Ilya Kerner, "Recofit: using a wearable sensor to find, recognize, and count repetitive exercises," in *Conference on Human Factors in Computing Systems*, 2014, pp. 3225–3234.
- [40] Huan Song, Jayaraman J. Thiagarajan, Prasanna Sattigeri, Karthikeyan Natesan Ramamurthy, and Andreas Spanias, "A deep learning approach to multiple kernel fusion," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2017, pp. 2292–2296.
- [41] Matthew D. Zeiler, "ADADELTA: an adaptive learning rate method," *CoRR*, 2012.
- [42] Yann LeCun, Léon Bottou, Genevieve B. Orr, and Klaus-Robert Müller, "Efficient backprop," in *Neural Networks: Tricks of the Trade - Second Edition*, pp. 9–48, 2012.
- [43] Panos P. Markopoulos, "Linear discriminant analysis with few training data," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2017, pp. 4626–4630.