

Impact of Featureset on Cross-Trials and Cross-Subjects Evaluation for Human Activity Recognition

Hosein Nourani ^{1,†,‡} , Emad Shihab ^{1,‡} and Omid Sarbishei ^{2,*}

¹ Dept. of Computer Science and Software Engineering; h.nourani@hotmail.com

² Dept. of Computer Science and Software Engineering; e.shihab@concordia.com

* Correspondence: e-mail@e-mail.com; Tel.: (optional; include country code; if there are multiple corresponding authors, add author initials) +xx-xxxx-xxx-xxxx (F.L.)

† Current address: Affiliation 3

‡ These authors contributed equally to this work.

Version July 25, 2019 submitted to Journal Not Specified

Abstract: Human Activity Recognition (HAR) refers to an emerging area of interest for medical, military, and security applications. However, the identification of the features to be used for activity classification and recognition is still an open point. In this work we have compared several state-of-the-art featuresets on HAR. To compare the featuresets, we have conducted an extensive set of experiments to indicate the vulnerable points in using featuresets. Aiming this goal, we implement several state-of-the-art machine learning models, ranging from traditional classifiers like SVM and GLM to convolutional neural networks; the models get evaluated by the performance gained using two challenging evaluation methods including cross-trials evaluation and cross-subjects evaluation, in addition to the convention method of 10-fold validation. The data from 55 gym exercises performing by thirteen subject have been recorded. To keep the realistic condition, they were asked to perform based on their own practice program. To capture the movements, two IMUs have been attached to the right wrist and right foot of the subjects. In total, 1300 features were extracted from the data. Our results showed that FNN achieved the highest performances using different featuresets. Against the prevailing wisdom, the statistical features gets outperformed by histogram feature set to deliver the best performance under all evaluation methods. *[Hosein: After all, this work provides the intuitions behind 100 hand-crafted features for beginner and researcher in HAR]*

Keywords: feature selection; featureset; Wearable; Motion sensor; Neural Network, Histogram, Human activity recognition

1. Introduction

Human Activity Recognition (HAR) using inertial sensors has been a target of a lot of researches in the last decade. One reason is because it opens a door to the thousands of useful applications in healthcare, physical fitness monitoring, elder care support and etc. Another reason is due to the pervasiveness of wearable sensors. They are small in size, cheap, and already embedded in almost all wearable gadgets. This is a promising point for HAR applications since no extra hardware setup-cost is required for users to start using them immediately. This is to say that the most challenging part in HAR is to producing an application to be able to efficiently use the device's resources and deliver the most accurate result in the form of classifying the activity type or counting it [1].

Basically, a HAR system using Inertial Measurement Units (IMUs) is composed of two basic components: (1) a data acquisition unit responsible for capturing human movements, (2) a processing unit responsible for recognizing the certain activities among movements of the subject [2]. In the first

component, the human movements is captured by different motion sensors such as accelerometer, gyroscope, and so on. These sensors can be located in an off-the-shelf device like a smart-phone for general HAR applications or specifically are accompanied by a storage and a processor, formed a System on Chip (SoC) for a certain purpose.

The captured movements as raw data transmits to the second component for processing operation. Within this component, there is a pattern recognition (HAR) model that classifies the input signals into certain classes of activities. This HAR model typically consists of three phases. First, there is a pre-processing operation that extracts informative features from raw signal. Second, a classifier is trained over extracted features. Third, an evaluation method to ensure that the classifier provides the required performance[3]. In other words, a HAR model should address these three aspects to be able to recognize an activity. The previous studies have achieved many improvements in each stages[4] however there are challenging aspects that they did not contemplate them and they are still open:

- Identifying the correct set of input features for the classifier[2] (in the phase 1)
- Existing an unique and reliable validation protocol to assess the models [4](in the phase 3)

While the first issue plays an important role to improve the performance of the model, the second issue is a critical point since different validation methods show the performance differently, somehow biased. For the same reason, using HAR models on new dataset shows a significant decrease in performance. As a consequence of this issue, currently it is impossible to know the state-of-the-art methods in human activity recognition[4]. Addition to this, some factors like temporal changes in user body (e.g., tiredness) have impacts on how they perform an identical activity within two different trials. However, it is impossible to evaluate such impacts by using the the traditional validation methods like k-fold cross-validation or Leave-One-Subject-Out (LOSO). To address this issue, in this study, we employ Leave-One-Trial-Out (LOTO) cross-validation addition to conventional validation methods. Using this validation method, we show that how the performance and robustness of the model change by feeding it by more trials.

Regarding the first issue, previous works have employed a wide range of features for HAR model. Dealing with selecting appropriate features for a HAR model (feature selection methods) are usually aiming to improve the performance of the model. However in [2] which is the most similar previous work to us, the authors also listed three extra issues which less have been considered in previous works: (1) The processing cost of producing a feature (2) The complexity of calculation of features makes a model difficult to understand. (3) The negative impacts of having too much number of features on performance of the model. In our previous work[5], we showed that by removing redundant features, we can keep the performance identical while shrink the dataset to 8% of its initial size. This approach significantly decreases the processing cost and vulnerability of getting overfit. In this study, we will focus on decreasing **the complexity of calculation of features** rather just decreasing size of dataset in a feature-set.

The aforementioned discussion motivated our study, where we evaluate a wide range of features by the performance reached by six popular classifiers. The models train and evaluate on two huge datasets (include our dataset which is one of the biggest dataset on gym exercises, publicly available). To provide a more robust evaluation we validate our models under several cross-validation methods.

The aim of this article is to present the following main contributions:

- Implementation and evaluation of the remarkable state-of-the-art featuresets.
- Analysing the impact of temporal changes on HAR model.
- Proposition of a novel model using Forward Neural Network (FNN) in HAR

Regarding the target datasets, we chose [] dataset since it provides a wide range of activities (33 activities) and it is big enough. However, to the best of our knowledge, there is no dataset that targets different trials of same activity with same subject, publicly available. So, we collected a dataset of 15 subjects and +50 gym activities including different trials for each activity.

The rest of this paper is organized as follows. Section 2 describes the study setup including the approach and the dataset used in this study. Section 3 defines the featuresets and their advantages in previous works. Section 4 presents the results of using featuresets on classifiers. Finally, Section 7 summarizes this study and highlights the most important contributions of this paper for future works.

2. Related Work

To the best of our knowledge, the statistical functions on time-domain signal provide the most popular features in HAR. Features like mean, median, mean standard deviation, variance, minimum, maximum are the most popular ones. The second group, Frequency-domain features are so helpful for those models that target periodic activities like walking. To achieve a frequency-domain feature, one segment of time signal should be transferred to several components in frequency domain. Each component may/may not have meaning depend on the activity. Frequency bins and auto-correlation coefficient are some examples.

Suto et al., 2017 investigate the efficiency of the popular machine learning strategies based on a right-ankle-mounted accelerometer, and their results suggest that one sensor is not enough for appropriate daily activity recognition due to the similar data generated from one sensor for different activities. Generally, One to One is the basic deployment and more suitable for the basic recognition tasks, such as step counting or sleep quality monitoring. Placing more sensors on multiple body parts is intuitively beneficial for improving the performance and robustness, whereas this can also result in increased complexity in deployment and computation cost. Sztyley et al., 2017 develop a position-aware HAR system by placing seven accelerometers in different body positions.

3. Data

The most crucial requirement to start an activity recognition process is having the real-world data of human activities. Although there are a couple of datasets publicly available for HAR[6], to the best of our knowledge, none of them provide the data of different sessions (neither at same day or different days) of same subject. In addition, they mostly focused on daily routine activities rather than transition activities (e.g., gym exercises) which are repetitive and more complicated in terms of pattern recognition[6]. Therefore, in this study we collect a dataset providing the following features:

- Specifically targeted on regular gym activities (including 55 exercises)
- The data tracks activities of a set of subjects over two to six weeks
- Data is recorded by four sensors (two wrists and two feet)
- Subjects performed the practices only based on their own experiences (there is no instruction)

To aim this goal, the first step is to determine the sensor type and how to deploy a system to collect the data. In this work, we have employed a SoC device which is called *Neblina*.

3.0.1. Neblina

Neblina is a miniature-sized box containing three tri-axial motion sensors (accelerometer, gyroscope, magnetometer) along with a processor, a flash memory, battery, and a bluetooth port. Using blue-tooth port, it can transmit the result to a host (e.g., cellphone or desktop computer). In fact, *Neblina* is equipped with all requirements for a real-time HAR system. Comparing with a smartphone, *Neblina* is much smaller (Figure 1) that lets us to attach it to different part of the subject's body without making any interrupt in his/her actions[7]. Having access directly to different resources like sensors or memory without OS interferences is another advantage of using *Neblina* that let us improve the efficiency of the model.

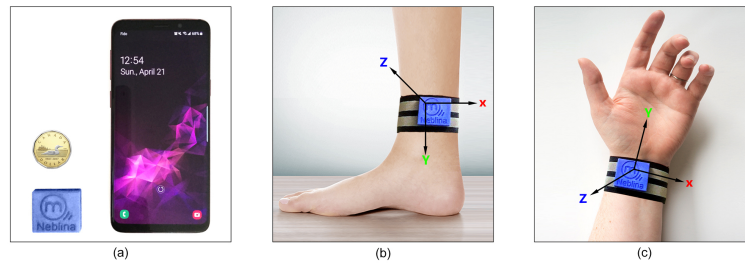


Figure 1. Neblina setup. (a) Compares dimensions of Neblina with a 1 dollar coin and a cellphone (Samsung Galaxy S9). (b) How Neblina located on foot using a strap. (c) How Neblina located on wrist using a strap.

3.0.2. Sensors

Depending on how much an activity is complicated (e.g., how many part of body are involved or how many stages are involved in it), a researcher may need to attach one or more sensors on different positions of the human body. However, using more sensors affects usability negatively. From the literature, using sensor on wrist for most of upper activities and using a sensor on foot for lower-body activity are more effective than other locations including chest, waist, thigh and so on. In this study, to cover activities both on lower-body and upper-body, two sensors were attached to right wrist and right foot (Figure 1 (b) and (c)). Although the device provides the magnetometer signals, we limited our process on using the accelerometer and gyroscope signals only. It is because the magnetometer signal can be affected by getting close to iron equipments in the gym. The frequency rate is fixed on 50Hz. It is worth to mention that the frequency rate more than 50Hz is not necessary because according to the Nyquist theorem, this rate is enough to record a repetitive activity with 25 cycle per second which is so much faster than the iterations of normal workouts in the gym (one iteration per 1-5 seconds).

3.1. Subjects

We asked 15 members of a gym (4 female), ages 21-35, to participate in this study. Participants varied in level of expertise (from 1 month to 6 consecutive years of experience). For more realistic scenarios, we did not constrain participants to certain exercises, instead we asked them to follow up their own plan. Comparing with previous works' dataset, considering this level of freedom for subjects returns following advantages: (1) Since each session is about 1 to 2 hours, we can observe the impact of fatigue on performing an activity. (2) The unknown period or null-class activities are not artificially performed, since subjects were free to do whatever they normally do in gym (3) impact of background experience can be measured. It is because the gym programs are cyclic over week or month. By repeating an activity over cycles, Subjects will be more consistent over different sets (a consecutive sequence of doing one activity). (4) there is no instruction of how doing exercises for subjects. Although this can let subjects to perform an activity in non-identical way, it is considered as an advantage for our study since it replicates the real-world condition. In [8], the authors showed that by changing the environment from space-constrained laboratory to a real gym the segmentation performance for recognizing gym exercises has dropped by 50%. Therefore, another advantage of keeping the experiment under real-world condition is the performance of the model is more reliable.

3.2. Activities

Our dataset ended up with 55 common exercises in gym. During electing participants, we picked mostly those persons who do more common exercises involved in either upper body or lower body. Thus, activities like Wall ball, jumping jacks and so on or advanced exercise in body-building are not in our dataset. In [9] (the second dataset in this work), authors have targeted CrossFit activities which are involved in upper and lower body together. They have shown that only one sensor on wrist is

enough to recognize such activities. Therefore, in this study we focused on those exercises in which either lower or upper body keeps stable during the activity. Thereby, existing the second sensor is necessary for recognizing the activity. As listed in Table 1, only two activities are involved in both upper and lower body. Running on treadmill (A1) make lower body involved. While In previous works this activity was recognized by sensor on wrist. Since a user can put her/his hands on device handler, using sensor on wrist is not effective always.

Table 1. List of exercises along with target body part involved in each exercise.

Exercise	Body Involved	Code
Treadmill	lower	A1
Ab crunch machine	lower & upper	A2
Lying leg curl	lower	A3
Barbell bicep curl	upper	A6
Standing calf raise	upper	A7
Seated calf raise	upper	A11
Overhead dumbbell press	upper	A12
Machine shoulder Press	upper	A13
Overhead barbell press	upper	A14

3.3. Data points

To generate the data points, previous works have employed different strategies. One well-known method for time-series signals is sliding-windows. As long as an activity is squashed in a range of samples during time, a model can see a stream of recorded data through a window with limited length of seconds (e.g., in this study it is 5 seconds time window). This window slides through the stream with a certain step size called shifting size. As long as the shifting size smaller than widow size, the sliding window is called overlapping and non-overlapping if they are equal. Previous works have shown that the different lengths of window size and shifting size influence the performance of the model and the computational cost. Because the activities in this study are gym exercise which the do not take longer than 5 seconds, intuitively, we choose 5 seconds for window size. It is a safe window size to ensure that at least one cycle of the activity can be completely seen in a window frame. we defined 200 milliseconds for shifting step which keeps the model more sensitive against changes in signal at the expense of more computational cost. Having such small shifting step does make sense since in real-world applications it decreases the latency of the application on predicting the activity type. Addition to the time period, the window length can be defined by the sampling rate of sensor. In this work, the sample rate is 50Hz. So, each window contains 250 ($5 * 50$) samples.

3.4. Dataset

To label the data we employed a process including three phases: (1) Before beginning of each session, each subject was asked to fill a form about list of activities, number of sets, and the weights if applicable. (2) During the session, a supervisor manually records type of exercise, the moment of start and stop, and number of repetitions. (3) After finishing the session, in order to have our desired accuracy in labelling, we visually trace the signals of accelerometer and gyroscope to refine the regains assigned to each exercise. Table 2 shows the statistics of the dataset. Since subjects may participate in more than one session, next column after subjects, shows the total number of sessions for each activity. , in the initial dataset, the number of subjects who are involved in all exercises is not equally distributed. Although Thus, we defined four datasets corresponding with our four experiments including K-fold evaluation, Leave-One-Set-Out evaluation, Leave-One-Session-Out evaluation, and Leave-One-Subject-Out evaluation.

Table 2. Statistics of the dataset divided by type of exercise along with the experiments that involve them in.

Exercise	Subjects	Sessions	Trials	Reps	Samples	10-fold	LOSO	LOTO
Treadmill	7	17	17	17	1M	Yes	Yes	Yes
Ab crunch machine	7	17	17	17	1M	Yes	Yes	Yes
Lying leg curl	7	17	17	17	1M	Yes	Yes	Yes
Barbell biceps curl	7	17	17	17	1M	Yes	Yes	Yes
Standing calf raise	7	17	17	17	1M	Yes	Yes	Yes
Seated calf raise	7	17	17	17	1M	Yes	Yes	Yes
Overhead dumbbell press	7	17	17	17	1M	Yes	Yes	Yes
Machine shoulder Press	7	17	17	17	1M	Yes	Yes	Yes
Overhead barbell press	7	17	17	17	1M	Yes	Yes	Yes
Non-activity	7	17	17	17	1M	Yes	Yes	Yes

4. Method

A typical work-flow on human activity recognition is including pre-processing, feature extraction, classification and evaluation. The process is shown in Figure ???. In this work we investigate the impact of 5 different featuresets in feature extraction phase; 5 different classifiers on classification phase; and 3 validation methods in evaluation phase.

4.1. Pre-processing

In this work, the common pre-processing operation for all experiments is a scale normalization. In fact, we scale all the input signals into a range between 0 and 1. Nevertheless, in some cases based on the requirements of the featuresets we add an extra layer of pre-processing which will be explained in related section.

4.2. Feature Extraction

In literature, in order to gain more information from the sensory data, the authors came up with hundreds of hand-crafted features among time domains, frequency domain or a combination of them. While each feature represents the signal in a certain point of view, it does not mean that this feature is informative enough for a model to recognize an activity based on that. One decisive factor to build a new feature is respecting the type of activity. Wang et. al. in [6] categorized human activities based on velocity and complexity (number of phases) of an activity into three main groups: (1) The **basic activities** which happen in comparatively longer duration e.g., walking and running. (2) The complex activities that are in the form of a sequence of several phases. Each phase might be a complex or basic activity e.g., coffee time, smoking. (3) Transition activities which having a certain but temporal pattern happening between two different postures or two basic activities. e.g., stand-to-sit, push-ups and so on. From this point of view, previous works introduced different feature sets that each one target a given type of activity. Consequently, targeting more than one type of activity brings more challenges for researcher to create a suitable set of features. Since exercises in gym composed of an orchestration of different type activities (basic, transition, complex), presenting a model to effectively recognize these activities will be more challenging. In this work, we compare the most state-of-the-arts feature sets in the literature. It is important to mention that we selected the ones that provide enough information to reproducibility (e.g., the definition of the features).

Table 3. Statical Functions along with the definitions and abbreviations

Code	Function	Description/Intuition	abbreviation
S1	Minimum	The value of the least sample	MIN
S2	Maximum	The value of the greatest sample	MAX
S3, SS8	Mean	The average of all samples	MEA
S4	Median	The middle value of samples	MEA
S5	Mean Absolute Deviation	The average distance between each sample and the mean of the stream	MAD
S6	Median Absolute Deviation	The average distance between each sample and the median of the stream	MAA
S7	Inner Quartile Range	The amount of spread in the middle part %50 of the stream	IQR
S8	Mean Crossing Rate	The rate of passing the mean along the stream	MCR
S9, SS9	Standard Deviation	how far the samples are from the mean	SD
S10, SS10	Variance	the average degree of distance between samples and mean	VAR
S11, SS11	Root Mean Square	The square root of the arithmetic mean of the squares of samples	RMS
HB	Histogram Bin	a 20 bins distribution of data	Hbin (1-20)
SS1	Number of autocorrelation peaks	The greater number of peaks refers to non-periodic activity and vice versa.	NACp
SS2	Prominent autocorrelation peaks	NACp with an extra condition that the peaks should be greater than neighbours with at least a certain distance	NACPP
SS3	Weak autocorrelation peaks	NACp with an extra condition that the distance between the peaks and neighbours should be less than a certain distance	NACWP
SS4	Maximum autocorrelation value	Value of the greatest peak (except for the initial peak at zero lag)	MAXAc
SS5	Height of the first autocorrelation peak (after zero-crossing)	less height refers to more fluctuations within the stream	FACp
SS6	Power bins (10 bins)	A 10 bins distribution of amplitudes of frequencies from 0.2-25Hz	Pbin(1-10)
SS7	Integrated RMS	The root-mean-square amplitude of the signal after cumulative summation	IRMS
Ph1	Movement Intensity	the Euclidean norm of the total acceleration vector after removing the static gravitational acceleration	MI
Ph2	Normalized Signal Magnitude Area	the acceleration magnitude summed over three axes within each window normalized by the window length	SMA
Ph3	Eigenvalues of Dominant Directions	The eigenvectors of the covariance matrix of the acceleration data correspond to the dominant directions along which intensive human motion occurs.	
Ph4	Correlation between Acceleration along Gravity and Heading Directions	It shows the human movement is either vertically or horizontally.	CAGH
Ph5	Averaged Velocity along Heading Direction	The Euclidean norm of the averaged velocities along y and z axes over the window.	AVH
Ph6	Averaged Velocity along Gravity Direction	averaging the instantaneous velocity along the gravity direction at each time t over the window	AVG
Ph7	Averaged Rotation Angles related to Gravity Direction	The cumulative rotation angles around gravity direction	ARATG
Ph8	Dominant Frequency	The frequency corresponding to the maximum of the squared discrete FFT component magnitudes of the signal from each sensor axis	DF
Ph9	Energy	The sum of the squared discrete FFT component magnitudes of the signal from each sensor axis	ENERGY
Ph10	Averaged Acceleration Energy	The mean value of the energy over three acceleration axes	AAE
Ph11	Averaged Rotation Energy	The mean value of the energy over three gyroscope axes.	ARE
OI1	Orientation Independent	result of applying PCA on Single Value Decomposition of x/y/z values of the stream	PCASVD(1-30)

4.2.1. Set_A: Statistical Features

Time-domain features are those features obtained directly from a window of sensor data and are typically statistical measures. They have been intensively investigated in different applications and proved to be effective and useful for HAR. These features are based on a comprehensive and intuitive understanding of how a specific activity or posture will produce a set of discriminative features from measured sensor signals. Lots of studies used statistical features in their work because of the reasons mentioned above as well as the low computational cost of them. Set A composed of 20 features based on applying 10 statistical functions on 12 input time-domain signal including (x/y/z axes of accelerometer and gyroscope, along with the cumulative sums of those axes). Table 3 shows the functions and a short description about each of them.

4.2.2. Set_B: Histogram bins Features

The histogram representation of a time series signal is equal to the probability distribution of that signal over a period of time (window size). In HAR, considering the fact that each activity contains a set of small movements (as small as one sample) with certain acceleration and rotation, histogram can indicate the difference between activities by showing the different distributions of those small movements. Comparing with other statistical features, extracting histogram bins is low cost in terms of required processing time and memory usage[10]. In this work we consider 20 bins to distribute the data among them (HB in Table 3)

4.2.3. Set_C: Self-Similar Features

Although by having enough number of bins, histogram can show the trend of data, they can not display the individual times, nor the iterations within the data. In addition, an exercise activity is inherently repetitive rather than a non-exercise activity. In other word, the signal from non-exercise activities looks more stochastic in a short period of time. Set C is designed based on the idea of extracting repetitions from the signal. This attribute of human activity can be extracted by calculating the convolution of a signal with shifted version of itself (autocorrelation) as well as by extracting the components of signal at frequency domain. Authors in [8] leveraged both methods (the autocorrelation and frequency functions) to make a featureset (self-similar features) to recognize gym exercises. This feature-set contains 160 features. They prepared 8 input signals including:

- The x axis of each sensor
- The magnitude of the accelerometer and gyroscope.
- The first principal component of all axes of each sensor.
- The first principal component of y and z axes of each sensors.

Table 3 shows the list of all features that are employed in this feature set (SS1-SS11) along with the intuition behind each of them.

4.2.4. Set_D: Physical Features

One idea to design more informative features from sensory data is to bring the principles of human movements into consideration. Specifically this type of features are derived based on the physical interpretations of human motion which called physical features. In [11], the authors introduced a set of features based on physical parameters of human motion. To have an strong physical meaning of motion data (e.g., moving forward, backward), they assumed that the sensor position and direction are known during the experiment. Comparing with previous featuresets, these features made up from a fusion of multiple inputs rather just applying a function on one input. As a part of pre-processing

operation for some of their features, they mentioned that they remove gravity from acceleration but they did not reveal any formula for that. In this paper, we remove the gravity using gyroscope data by applying the method described in [12].

4.2.5. Set_E: Orientation Independent Features

In contrast to physical features which is dependant to the position and orientation of the sensors, in [13], Yurtman et. al. targeted on a HAR model that does not rely on the variation of sensor orientation. In fact, in their model, they introduced an Orientation-invariant transformations (OITs). They compared their model with the ordinary model - pre-defined sensor orientation, on five different datasets. Although their featureset did not have a significant impact on performance, it brought an extra added value to the model that lets the model to be more robust against orientation. The OIT that they have introduced in their work is inspired by the idea of *single value decomposition*[14]. First, they project every data point from original x/y/z space to a new space with same number of dimensions but with largest distance between data points. The intuition here is that the direction of the axes are defined by value of the data points not by x, y or z direction. Next, they apply PCA on transformed data and take first 30 most informative features. In Table3, these type of features are indicated by "OI" prefix.

4.3. Activity Recognition

[Hosein: high-level Description => Based on Activity Recognition Chain (ARC) in [15]] [Hosein: Comparing the process between classical models VS Neural Network model]

In the well-known paper [15], the authors have described the main characteristics of a HAR model into five categories. They categorized the HAR models by: (1) execution type (offline/ online), (2) type of the activities that the model can recognize (Periodic/Sporadic/Static), (3) type of input signal (segmented/ continuous) (4) dependency of model to the user (user-independent/ user-specific), (5) dependency of model to the other inputs like user's context (stateless/ stateful). In this study, models are stateless and work on segmented stream in offline mode. Regarding dependency of model to the user, in this study, we investigate different aspects of user that can influence the performance of the model.

4.3.1. Support Vector Machine

Assuming each data point is a co-ordinate (support vector) of feature space, the Support Vector Machine (SVM) is an algorithm to find an optimum hyperplane between support vectors between two classes. To create a more robust classification, SVM maximises the distance between the hyperplane and the closest support vectors of each class (margin); This happens at the expense of increasing the tolerance(mis-classification) of the model (loss function). In a situation that finding a hyperplane between two classes is impossible, the SVM project the feature space to another space having more number of dimensions. This is the responsibility of a function called Kernel function.

4.3.2. Decision Tree

Assuming a training set size of N samples and M features, an ensemble of K decision trees (weak learners) is called RF such that in the training phase, each tree in the classifier randomly picks a set of training samples ($n \ll N$) via bootstrapping and it finds the best split from a set of randomly selected features ($m \ll M$). To get the class labels in the test phase, the class probabilities at leaf nodes of the trees are averaged and the highest average probability is chosen as the class label of the input sample.

4.3.3. K-nearest neighbour

KNN algorithm is a simple classification algorithm based on the calculation of the distance (usually the Euclidean distance) between the new element to be classified and the elements in the training set. Firstly, the training elements are sorted in descending order according to their distance from the new element. Then, the most frequent class of the first K elements (called neighbors) is associated to the new element. For this kind of classifier, only the value of the K neighbors must be decided. A common starting N, where N was the number of elements in the training set. Beginning from value for K is K_{in} this consideration, we decided to analyze 32 values around K_{in} and, thus, we used five bits for the second substring of each GA solution ($2^5 = 32$): each possible value assumed by the second substring was associated to a specific K value to be set in the classifier.

4.3.4. Forward Neural Network

A FNN is made up of a set of neurons, connected by weighted arcs, that process the input information according to the McCulloch and Pitts model:

$$y = f\left(\sum_i w_i \cdot x_i\right) \quad (1)$$

where y is the output of the neuron, w_i are weights of the incoming connections, x_i are inputs to the neuron, and f is called transfer function and should be selected according to the classification problem. Neurons in a FNN are organized in layers: in the input layer, one neuron for each input variable is required; the number of neurons in the output layer is decided according to the number of classes to be recognized and the selected transfer function; between input and output layers a certain number of hidden layers can be inserted, whose dimensions are usually decided testing different configurations. In this study we fixed a basic network structure with input layer and first hidden layer both including one neuron for each feature selected according to the first substring of the GA solution, and an output layer made up of one neuron returning the recognized activity. Then, the number of hidden layers was increased according to the second substring of each solution: three bits were used for adding from one to eight further hidden layers to the basic structure. Each new hidden layer included 1/2 of the previous layer neurons. The sigmoid transfer function was used for all hidden layers and the linear transfer function was set for the output neuron. Since the output neuron returned a real value for each classified element, the round operator was applied to the FNN output and used to assign the final class.

4.4. Evaluation

There exists a set of metrics that are more popular to measure the activity recognition performance, such as accuracy, recall, precision, F-measure, and so on. Assuming the basic terms in confusion matrix have been defined as follows:

True Positive (FP): These are cases in which we predict an activity, and user was doing that activity.

True Negative (TN): In which we predict a non-activity period, and user was not doing a certain activity.

False Positive (FP): In which we predict a certain activity for a segment of data, however, user is either doing another specific activity or generally doing something else (out of activity given list)

False Negative (FN): In which we predict either a not-activity period or a certain activity, but, it is not the activity that user is really doing that.

The most popular metrics are the following:

Accuracy measures how often the classifier is correct. Specifically, it is equal to $(TN + TP) / \text{total}$.

Miss-classification measure how often the classifier is incorrect. Specifically, it is equal to $(1 - \text{Accuracy})$.

Precision measures when the classifier detects an activity, how often it is correct. Specifically, it is equal to $TP / (TP + FP)$.

Recall measures when user is doing a certain activity, how often the classifier can detect it correctly. Specifically, it is equal to $TP / (TP + FN)$. This term is also known as *Sensitivity* or *True Positive Rate*.

F-Score measures a weighted average of both Recall and Precision. Specifically, it is equal to $(2 \times \text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$.

Null Error Rate: measures how often it is incorrect if we constantly return the major class in dataset as response of the classifier. To the best of our knowledge, in most HAR datasets, the major class is non-activity class.

Table 4 summarizes the main evaluation metrics employed in some recent works. Accuracy is the obvious alternative shown in the table. F-measure is another popular metric followed by Precision and Recall with less popularity. Although accuracy is a straight forward measure to show the performance of the model, in case of having an unbalanced class within the dataset - which is often the case on wearable applications, it may get biased easier than other metrics. In this situation, comparing Null Error Rate for the unbalanced class with the accuracy of that class can explain the real performance of the HAR model. On the other hand, F1 is a suitable measure since it relies on both the precision and recall; knowing, the recall is less affected by unbalanced classes.

Table 4. Evaluation metrics and methods in previous works

Work	Metrics	Method
[16]	Accuracy	4-fold
[5]	Accuracy, F1	10-fold
[17],[9]	Accuracy	5-fold
[8]	Precision, Recall	Leave-One-Subject-Out
[11],[18]	Accuracy	Leave-One-Subject-Out
[13]	Accuracy	Leave-One-Subject-Out, 10-fold
[2]	Accuracy, F1	Student t-test

After choosing an appropriate measure, we need a validation method to asset the model by that measure. Aiming this, we separate the available data into training and testing sets. For this purpose, in the context of HAR, there are traditional techniques such as k-fold cross-validation, leave-one-subject-out as well as relatively less common techniques like hold-out and leave-one-trial-out available in the literature. In this study, we pick three validation methods which almost cover main ideas in evaluation of a HAR model in recent works. The other methods like leave-one-sample-out or holdout can be considered as different forms of these three methods.

4.4.1. K-fold Cross validation

The most typical approach to evaluate the performance of HAR model is k-fold-cross validation. The idea is to use resampling procedure in a way that all the samples be used once during the testing period, mostly in the case of having limited dataset. The so called parameter k refers to the number of groups that the given dataset is to be split into. Each time one group becomes the test set and remaining k-1 groups become training set. During k turns, we evaluate the model k times on different test set and train set. Finally, the performance is summarized by averaging the performance of all k turns. In this work, we apply k-fold cross validation on all the models

4.4.2. Leave-One-Subject-Out Cross validation

We highlight that the Leave-One-Subject-Out (LOSO) protocol 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, the LOSO protocol reflects a realistic scenario, where a model is trained in an offline way, 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 performance from a subject to another, since the same activity can be performed in different ways by the subjects.

4.4.3. Leave-One-Trial-Out Cross validation

The Leave-One-Trial-Out (LOTO) cross validation process is exactly same as LOSO. However, instead of considering the subjects as folds, each trial (session) of doing the activity is considered as a fold. In other words, the recorded data from each session has an identifier to be distinguished from other sessions. So, during the evaluation, the dataset split into sessions instead of subjects. sessions can be for same subject or different subjects. In our dataset we also put an extra column showing the session id which is unique among all sessions of a given subject. The main advantage of using this process comparing with the previous one is that it needs less number subjects since each subject can have several sessions. And same as previous one, there is no overlap between train set and test set.

5. Results

5.0.1. RQ1: Which featureset delivers higher performance over classifiers?

In this section we compare the performance of four different classifiers on all five featuresets. Table 5 and 6 show the results respectively in Accuracy and F-score. We used 10-fold cross validation to measure the performance. A total of 150,000 data points were included in the validation set. Table 2 shows total number of samples for each activity separately. It can be seen how classifiers performance varies over featuresets. All classifiers on *set_B* and *set_A* deliver their best performance. On the other hand, the performance achieved on a given featureset has been affected by changing classifiers. FNN and SVM gained the best results comparing with two other classifiers. Also, it seems clear that the *set_E* (Orientation Invariant features) provide insufficient information for successfully tackling activity recognition, since the accuracy never exceeds 62%, whereas other featuresets achieve accuracies over 75%. Interestingly, the *set_B* (Histogram Features) with having the least complexity in computation gained the best performance among other featuresets. In fact, accuracy achieved by the top-performing models (FNN and SVM) are significantly better when they use histogram features. This could happen due to choosing best bin width based on length of activity and windows size.

Table 5. Accuracy for each classifier and for all each feature-set using 10-fold cross validation

Classifier	Set_A	Set_B	Set_C	Set_D	Set_E
SVM	78.45%	83.11%	74.16%	70.02%	59.38%
KNN	62.35%	81.11%	76.52%	59.72%	52.38%
FNN	86.45%	89.00%	72.51%	69.31%	61.88%
DT	69.94%	72.28%	73.76%	69.72%	61.38%

Table 6. F1 for each classifier and for all each feature-set using 10-fold cross validation

Classifier	Set_A	Set_B	Set_C	Set_D	Set_E
SVM	72.33%	79.08%	67.14%	66.45%	60.89%
KNN	60.17%	76.00%	73.23%	57.12%	53.98%
FNN	80.71%	84.38%	66.10%	66.13%	60.02%
DT	63.15%	67.99%	71.18%	63.61%	60.33%

5.0.2. RQ1: Which classifier performs better on gym exercise recognition?

5.0.3. RQ1: How are different evaluation methods?

6. Discussion

6.1. *K-fold Evaluation result is biased?*

As we argued earlier, each process has a drawback that might cause a negative impact on the methods. For instance, SNOW can produce biased results and the FNOW generates few samples. To face these problems, we propose the Leave-One-Trial-Out (LOTO) sample generation process. From the literature, the most popular method to evaluate a HAR model is K-fold cross-validation. In this method the data splits into k equal parts. The $k-1$ parts go for training and one part left for test. This process repeats for each k , separately. Main aim of this method is to keep the test and train part separate from each other and use all the samples in train and test. However, in case of using the conventional sliding window (with $n\%$ overlap) for generating data-samples it is impossible for samples to be completely separated from each other. It is because of the overlap period which is common between each two windows in a sequence. In other words, as mentioned in [4], the result is always biased because $n\%$ of the data are identical between test and train part. To avoid this drawback, a ordinary evaluation method is Leave-One-Subject-Out (LOSO). Basically, this method resembles the k -fold cross-validation in which a fold is replaced by a subject. This method is secure against being biased since the data from each subject has no common area with other subjects. However, in order to have the performance in a satisfactory level, we have to increase the number of subjects, respectively. The Leave-One-Trial-Out (LOTO) cross-validation[4] is similar to LOSO, but use the data of a trial instead of a subject. Thus, we have ensured that the samples are basically separated and the performance does not depend on the number of subjects any more. *[Hosein: + showing it in a example]*

6.2. *Hand-Crafted Features VS Automated Features*

Automatic feature learning

The extraction of hand-crafted features depends on domain knowledge. However, automatically learned features by the deep networks features are easy to understand and implement

The key advantage of using hand-crafted features is that the features are computationally lightweight to implement especially in ubiquitous devices (Morales et al., 2017). The strengths of the automatically learned features by the deep networks are that the learning can be very deep, and the learning process does not rely on domain knowledge.

Feature type Advantages Hand-crafted Features Automatically learned features Table 5 Comparison of hand-crafted features and automatically learned features Disadvantages Easy to understand the physical meanings of the features; Extraction is efficient and easy to deploy; Work well for many HAR problems. No domain knowledge needed; Automatically learning features from raw data; Features are more robust and generalized. Domain knowledge needed; Sensor-type specific; Need further feature selection. Lots of computing resources; Parameters are difficult to adjust; The learned features are less interpretable

Authors should discuss the results and how they can be interpreted in perspective of previous studies and of the working hypotheses. The findings and their implications should be discussed in the broadest context possible. Future research directions may also be highlighted.

Another way of evaluating the featuresets is using a feature selection method and compare the percentage of contribution of each featureset within the featureset result.

Using neural network in pattern recognition is one of the most promising topic in recent works, especially in activity recognition. It is because, the characteristics like flexibility of the model with more various activity sets or improving the performance by adding more data in training phase make

NN an appealing choice to data analyst. Related works have presented successful approaches on customizing neural nets to solve certain problems in HAR. In most studies, the authors introduced a new design of layers and sort of new configuration of neurons in each layer. Their results show that NN works better rather traditional models like GLM or SVM on more complex activities, in practice.

6.3. FNN training speed over different feature sets

Histogram bins progress outperforms all other featuresets in delivering better accuracy in less training time.

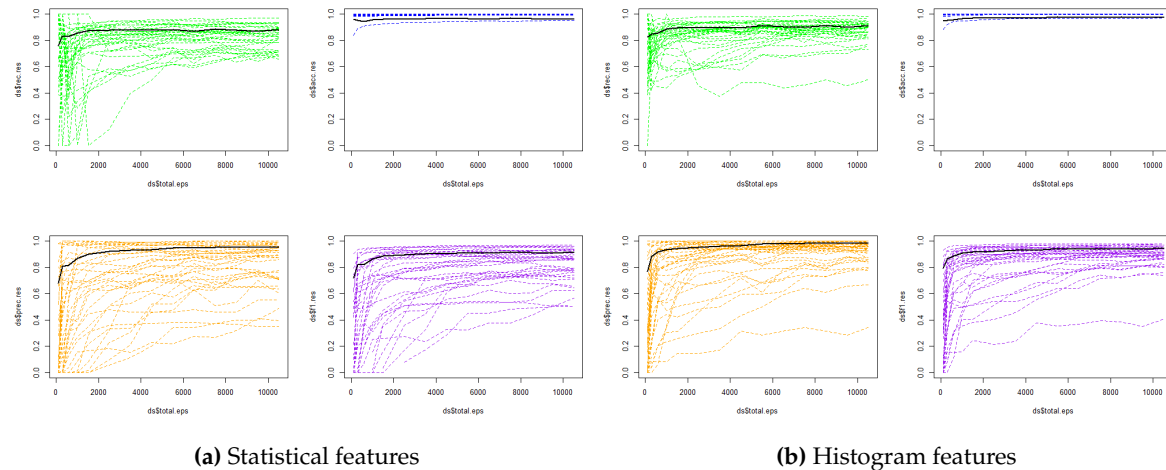


Figure 2. Comparing impact of featureset on FNN convergence speed

6.4. Impact of Imbalanced dataset

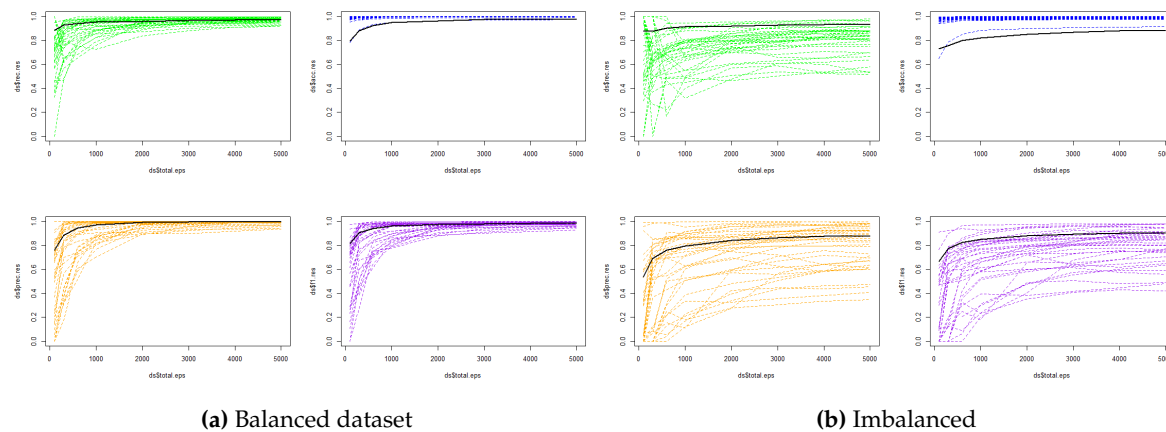


Figure 3. Comparing impact of undersampling on FNN

7. Conclusions

This section is not mandatory, but can be added to the manuscript if the discussion is unusually long or complex.

Author Contributions: For research articles with several authors, a short paragraph specifying their individual contributions must be provided. The following statements should be used “conceptualization, X.X. and Y.Y.; methodology, X.X.; software, X.X.; validation, X.X., Y.Y. and Z.Z.; formal analysis, X.X.; investigation, X.X.; resources, X.X.; data curation, X.X.; writing—original draft preparation, X.X.; writing—review and editing, X.X.; visualization, X.X.; supervision, X.X.; project administration, X.X.; funding acquisition, Y.Y.”, please turn to the [CRediT taxonomy](#) for the term explanation. Authorship must be limited to those who have contributed substantially to the work reported.

Funding: Please add: “This research received no external funding” or “This research was funded by NAME OF FUNDER grant number XXX.” and “The APC was funded by XXX”. Check carefully that the details given are accurate and use the standard spelling of funding agency names at <https://search.crossref.org/funding>, any errors may affect your future funding.

Acknowledgments: In this section you can acknowledge any support given which is not covered by the author contribution or funding sections. This may include administrative and technical support, or donations in kind (e.g., materials used for experiments).

Conflicts of Interest: Declare conflicts of interest or state “The authors declare no conflict of interest.” Authors must identify and declare any personal circumstances or interest that may be perceived as inappropriately influencing the representation or interpretation of reported research results. Any role of the funders in the design of the study; in the collection, analyses or interpretation of data; in the writing of the manuscript, or in the decision to publish the results must be declared in this section. If there is no role, please state “The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results”.

Abbreviations

The following abbreviations are used in this manuscript:

MDPI	Multidisciplinary Digital Publishing Institute
DOAJ	Directory of open access journals
TLA	Three letter acronym
LD	linear dichroism

Appendix A

Appendix A.1

The appendix is an optional section that can contain details and data supplemental to the main text. For example, explanations of experimental details that would disrupt the flow of the main text, but nonetheless remain crucial to understanding and reproducing the research shown; figures of replicates for experiments of which representative data is shown in the main text can be added here if brief, or as Supplementary data. Mathematical proofs of results not central to the paper can be added as an appendix.

Appendix B

All appendix sections must be cited in the main text. In the appendixes, Figures, Tables, etc. should be labeled starting with ‘A’, e.g., Figure A1, Figure A2, etc.

References

- Schilit, B.N.; Adams, N.; Want, R.; others. *Context-aware computing applications*; Xerox Corporation, Palo Alto Research Center, 1994.
- Rosati, S.; Balestra, G.; Knaflitz, M. Comparison of Different Sets of Features for Human Activity Recognition by Wearable Sensors. *Sensors* **2018**, *18*, 4189.
- Kołodziej, M.; Majkowski, A.; Tarnowski, P.; Rak, R.J.; Gebert, D.; Sawicki, D. Registration and Analysis of Acceleration Data to Recognize Physical Activity. *Journal of Healthcare Engineering* **2019**, 2019.
- Jordao, A.; Nazare Jr, A.C.; Sena, J.; Schwartz, W.R. Human activity recognition based on wearable sensor data: A standardization of the state-of-the-art. *arXiv preprint arXiv:1806.05226* **2018**.
- Nourani, H.; Shihab, E.; Sarbishe, O. The Impact of Data Reduction on Wearable-Based Human Activity Recognition. *Proceedings of the 15th Workshop on Context Modeling and Recognition*. IEEE, 2019, CoMoRea '19, pp. 89–94.
- Wang, Y.; Cang, S.; Yu, H. A survey on wearable sensor modality centred human activity recognition in health care. *Expert Systems with Applications* **2019**.
- de Faria, I.L.; Vieira, V. A Comparative Study on Fitness Activity Recognition. *Proceedings of the 24th Brazilian Symposium on Multimedia and the Web*. ACM, 2018, pp. 327–330.

8. Morris, D.; Saponas, T.S.; Guillory, A.; Kelner, I. RecoFit: using a wearable sensor to find, recognize, and count repetitive exercises. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 2014, pp. 3225–3234.
9. Soro, A.; Brunner, G.; Tanner, S.; Wattenhofer, R. Recognition and Repetition Counting for Complex Physical Exercises with Deep Learning. *Sensors* **2019**, *19*, 714.
10. Sarbishei, O. A Platform and Methodology Enabling Real-Time Motion Pattern Recognition on Low-Power Smart Devices. *2019 IEEE World Forum on Internet of Things* **2019**, pp. 257–260.
11. Zhang, M.; Sawchuk, A.A. A feature selection-based framework for human activity recognition using wearable multimodal sensors. *Proceedings of the 6th International Conference on Body Area Networks*. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), 2011, pp. 92–98.
12. Accelerometer. <https://www.w3.org/TR/accelerometer/>. (Accessed on 07/03/2019).
13. Yurtman, A.; Barshan, B. Activity recognition invariant to sensor orientation with wearable motion sensors. *Sensors* **2017**, *17*, 1838.
14. Moon, T.K.; Stirling, W.C. *Mathematical methods and algorithms for signal processing*; Vol. 1, Prentice hall Upper Saddle River, NJ, 2000.
15. Bulling, A.; Blanke, U.; Schiele, B. A tutorial on human activity recognition using body-worn inertial sensors. *ACM Computing Surveys (CSUR)* **2014**, *46*, 33.
16. Brownlee, J. A gentle introduction to k-fold cross-validation. *Accessed October* **2018**, *7*, 2018.
17. Shakya, S.R.; Zhang, C.; Zhou, Z. Comparative Study of Machine Learning and Deep Learning Architecture for Human Activity Recognition Using Accelerometer Data. *Int. J. Mach. Learn. Comput* **2018**, *8*, 577–582.
18. Mehrang, S.; Pietila, J.; Tolonen, J.; Helander, E.; Jimison, H.; Pavel, M.; Korhonen, I. Human activity recognition using a single optical heart rate monitoring wristband equipped with triaxial accelerometer. In *EMBECE & NBC 2017*; Springer, 2017; pp. 587–590.

Sample Availability: Samples of the compounds are available from the authors.

© 2019 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).