

# ML Project : Machine Learning Methods Applied to Human Physical Activity Classification Using On-Body Sensors

Mohamed-El-Amine Seddik

Marwen Sallem

Kais Slimi

Oussama Bouraoui

February 1, 2017

## Abstract

This paper describes our choice for the machine learning project subject, we chose to work on applying machine learning methods to the problem of human physical activity classification, based on data that comes from on-body sensors. Basically, the aim of this project is to explore different techniques of machine learning and to apply them to this particular problem.

## Motivation

The ability of designing a system that automatically classifies the physical activity performed by a given person is very attractive for many applications in various fields such that health-care monitoring and in developing smart human-machine interfaces. There are several techniques to handle this task, and mainly some of them are based on computer vision where some others aims to use on-body sensors. Based on [1] we chose to work on data that comes from on-body sensors. In fact, there is many datasets available on the internet that are used for this special problem, we provide below a link to a dataset that we plan to use for our project.

## Goal

As illustrated by the title, our main goal is to classify the Human Physical Activity. To reach this goal we plan to explore different machine learning methods (Probabilistic based methods like HMM and geometric based ones as SVM k-NN ...) and to apply them to this particular task, at the end we aim to provide the performances of the different applied methods.

## Dataset

- <https://archive.ics.uci.edu/ml/datasets/Activity+Recognition+from+Single+Chest-Mounted+Accelerometer>

# Introduction

Several techniques have been used to recognize Human activities such as video-based sensors, wearable-based sensors, environmental sensors and object sensors like smart phones. In this paper, wearable-sensors based systems for activity recognition have been used.

Human activity recognition or (HAR), begins by collecting the raw data from sensors worn on certain parts of the body. The sensors will provide different data depending on the part of the body where it's worn (for example wrist, chest, ...) and the technical characteristics of the sensor itself. One of the biggest challenges of HAR is to overcome this variability in sensor characteristics. The sensors usually have limited processing power which makes it difficult to get a high frequency sampling or do any kind of preprocessing on the sensors.[2]

We are therefore unable to use the raw data directly for the classification task. We need to do several preprocessing steps depending on the classification method we want to use. In section 1 we will explain the state of the art feature extraction techniques encountered in such problems.

Another issue in HAR is Intra-class variability: Different people will do the same activity differently. Their data on the sensors will not be very similar but it will have the same label (in the supervised setting) or the classifier has to assign them the same label (in the unsupervised setting). This can be overcome by collecting more data or develop person-independent features that are more robust to this variability. In the section 2, we will present a supervised technique that aims to classify the human activities using LDA and KDA classifier. But generally the supervised approaches require a large amount of labelled data. They are too expensive. For that, one has to rely on unsupervised approaches. Hence, we will present in 4 an unsupervised technique based upon joint segmentation of multidimensional time series using a Hidden Markov Model (HMM) in a multiple regression context and is denoted by MHMMR: Multiple Hidden Markov Model Regression.

## 1 Preprocessing and features extraction

Accelerometers provide three axis time series  $A_x$ ,  $A_y$  and  $A_z$ . Once sampling on a known frequency, we obtain a different signal for each class, for instance the figure 1

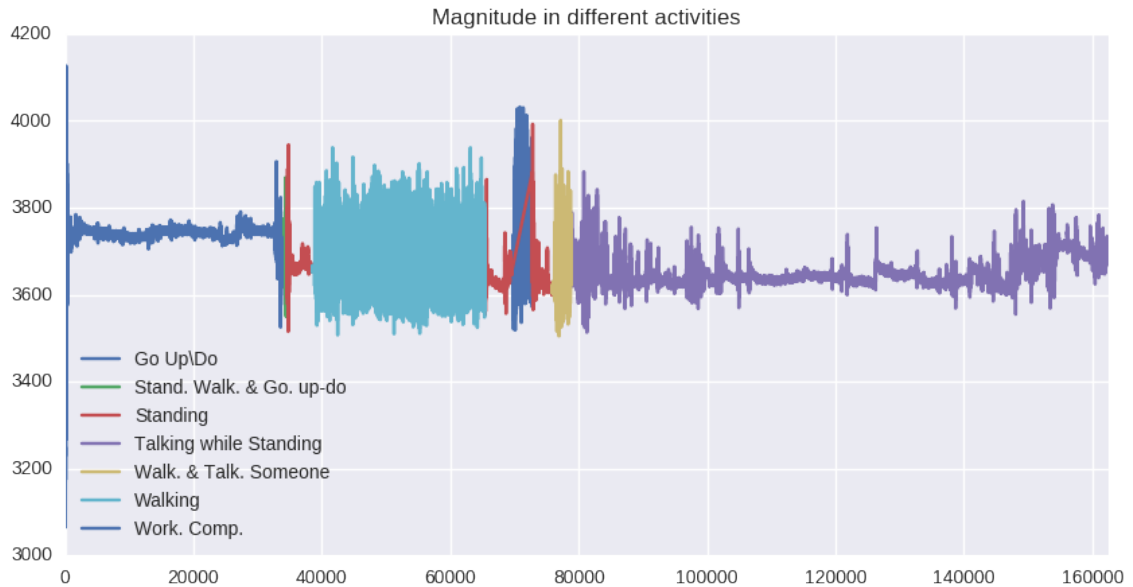


Figure 1: Accelerometer Data for Five Different Activities [1]

To be able to apply Machine Learning techniques, we have to extract features from these given signals. These latter is mainly done on two different phases :

In the first phase we are going to widen our signal databases by adding new signals :

- Compute the magnitude of the acceleration :  $A_m = \sqrt{A_x^2 + A_y^2 + A_z^2}$  [3].
- Separate every axis time-series in high-frequency (AC) components that captures dynamic motion, and low-frequency (DC) components related to the gravity's influence that captures static postures. Thereby, we will obtain for each time series three additional time series that corresponds to : the signal without filtering, the signal after applying a high-pass filter and the signal after applying a low-pass filter.

In fact, the magnitude figure 1

Now we have 12 signals extracted from the three initial signals  $A_x$ ,  $A_y$  and  $A_z$ . Having these new set of signals, we will split them into windows in order to compute features. In fact, we apply a windowing with overlapping process which means that we will take windows of 52 sample (1 second of data) with 50% of overlapping between windows. Then, we can move to the second phase, i.e extracting features :

- We compute the RMS of the velocity (integration of acceleration) and mean value of Minmax sums<sup>1</sup>. In [3] we showed that theses features are the most important for the proposed Random Forest classifier in the Table ??.
- Basic statistics as mean value, standard deviation, skewness, kurtosis, correlation between each pairwise of accelerometer axis. In fact, Manini and Sabatini shows in [1] that these features proved to be useful in HAR. The figure 2 shows that the basic statistics are good features to discriminate some classes.

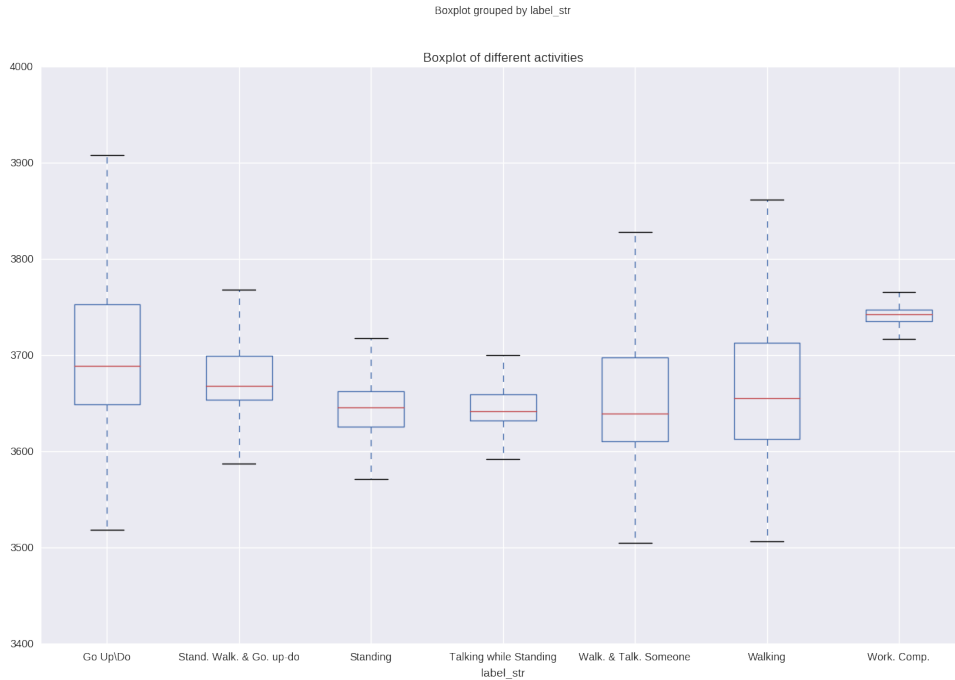


Figure 2: Means and variances for different classes.

- Other features like ARIMA coefficients [1] and energy of the coefficients of some level wavelet decomposition [3] are also of great importance.

Now, we will have more than 300 features. It's possible to use them either in Single-Frame approaches as " Linear and Kernel Discriminant Analysis approaches " or in Sequential approaches like " Hidden Markov Model Regression " which we will describe in the next sections.

<sup>1</sup>The Minmax sums are computed as the sum of all the differences of the ordered pairs of the peaks of the time series.[3]

**Table 1.** List of Features selected by Random Forest

Feature	Importance	Feature	Importance
Mean Value $A_{zdc}$	4.64	Mean Value $A_{ydc}$	3.86
MinMax $A_{zdc}$	4.61	Rms Velocity $A_{ydc}$	3.67
RMS Velocity $A_{zdc}$	4.23	Mean Value $A_{zb}$	3.59
RMS Velocity $A_{mdc}$	4.2	Mean Value $A_{xdc}$	3.57
RMS Velocity $A_{xac}$	4.14	MinMax $A_{xdc}$	3.52
Mean Value $A_{mdc}$	4.07	MinMax $A_{zb}$	3.51
MinMax $A_{ydc}$	3.92	Mean Value $A_{yb}$	3.33
Standard Deviation $A_{xb}$	3.9	Rms Velocity $A_{xdc}$	3.22
MinMax $A_{mdc}$	3.89	Rms Velocity $A_{zb}$	3.2
Standard Deviation $A_{xdc}$	3.87	MinMax $A_{yb}$	2.96

## 2 Linear and Kernel Discriminant Analysis approaches

One observes the raw data of the chosen dataset, we acknowledge an overlapping between the different classes (see Fig 3(a)), this overlapping is due to the high with-in and low between-class variances. Based on [4], in order to fix this overlapping problem, we present in this section the LDA and KDA models, which are a supervised classification approaches that utilizes the class specific information maximizing the ratio of the within and between class scatter information. In particular, the KDA is a generalization of LDA applied in a RKHS feature space.

### 2.1 Theory about Linear Discriminant Analysis

#### 2.1.1 Two-classes problem

The LDA responds to the question of how do we use the label information in finding informative projections of our data? To that purpose it considers maximizing the following objective:

$$(\mathcal{P}) : \max_{\omega} J(\omega) = \max_{\omega} \frac{\omega^T S_B \omega}{\omega^T S_W \omega} \quad (1)$$

Where  $S_B$  is the between classes scatter matrix and  $S_W$  is the within classes scatter matrix. And are given by:

$$S_B = \sum_c (\mu_c - \bar{x})(\mu_c - \bar{x})^T \quad S_W = \sum_c \sum_{i \in c} (x_i - \mu_c)(x_i - \mu_c)^T \quad (2)$$

Where  $\bar{x}$  overall mean of the data.

$J$  being invariant with respect to re-scaling of the vector  $\omega \rightarrow \alpha\omega$ . Thus we can always choose  $\omega$  such that  $\omega^T S_W \omega = 1$ , so that the problem  $(\mathcal{P})$  could be written as:

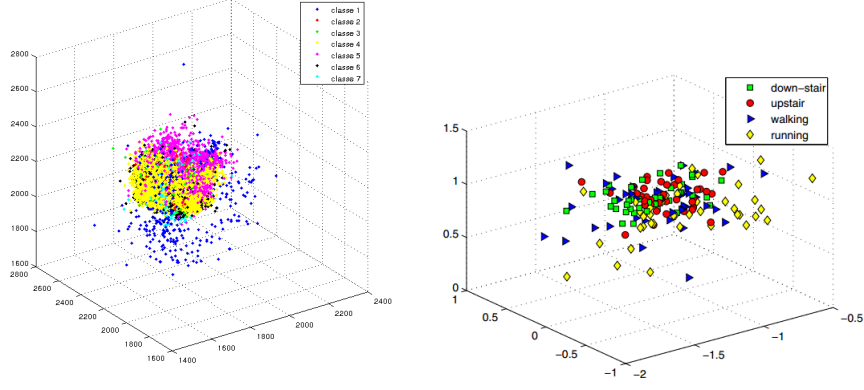
$$(\mathcal{P}) : \min_{\omega} -\frac{1}{2} \omega^T S_B \omega \quad \text{s.t.} \quad \omega^T S_W \omega = 1 \quad (3)$$

The Laplacian of  $(\mathcal{P})$  is thus:

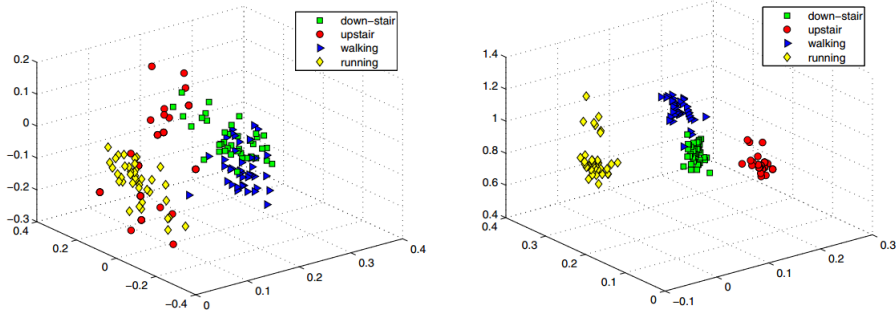
$$\mathcal{L}_{\mathcal{P}}(\omega, \lambda) = -\frac{1}{2} \omega^T S_B \omega + \frac{1}{2} \lambda (\omega^T S_W \omega - 1) \quad (4)$$

Finally, taking the derivative with respect to  $\omega$  of the Laplacian, one get:

$$S_B \omega = \lambda S_W \omega \Rightarrow S_W^{-1} S_B \omega = \lambda \omega \quad (5)$$



(a) Raw data plotted for different classes.



(b) Results of applying LDA (left) and KDA (right) to the raw data.

Figure 3: Visualization of the different classes before and after applying Linear and Kernel Discriminant Analysis, we notice that the classes are perfectly separable ones applying the KDA (see 3(b) right).

$S_B$  is symmetric positive definite (p.d.), thus  $S_B = S_B^{1/2} S_B^{1/2}$  where  $S_B^{1/2} = U \Lambda^{1/2} U^T$  (where  $U \Lambda U^T$  is the eigen-decomposition of  $S_B$ ). So performing the change of variable  $v = S_B^{1/2} \omega$ , one get:

$$\underbrace{S_B^{1/2} S_W^{-1} S_B^{1/2}}_{\text{symmetric p.d.}} v = \lambda v \quad (6)$$

Finally, the goal of the LDA is to find the eigenvector corresponding to the largest eigenvalue of the symmetric p.d. matrix  $S_B^{1/2} S_W^{-1} S_B^{1/2}$ . Denote this vector  $v^*$ , thus  $w^* = S_B^{-1/2} v^*$ .

### 2.1.2 MultiClasses Problem

Based on two classes problem, we can see that the LDA generalizes grace-fully for multiple classes problem. Assume we still have a set of  $d$ -dimensional samples  $X = \{x^{(i)}\}_{i=1}^n$ , and there are totally  $C$  classes. Instead of one projection  $y$ , mentioned above, we now will seek  $(C - 1)$  projections  $[y_1, y_2, \dots, y_{C-1}]$  by means of  $(C - 1)$  projection vectors  $\Theta$  i arranged by columns into a projection matrix  $\Theta = [\theta_1, \theta_2, \dots, \theta_{C-1}]$ , where:

$$y_i = \theta_i^T X \Rightarrow y = \Theta^T X \quad (7)$$

As in the two classes problem, the optimal projection matrix  $\Theta^*$  is the one whose columns are the eigenvectors corresponding to the largest eigenvalues of the following generalized eigenvalue problem:

$$\Theta^* = [\theta_1^*, \theta_2^*, \dots, \theta_{C-1}^*] = \arg \min_{\Theta} \frac{|\Theta^T S_B \Theta|}{|\Theta^T S_W \Theta|} \quad (8)$$

$$\Rightarrow (S_B - \lambda_i S_W) \Theta_i^* = 0, \quad i \in \{1, \dots, C-1\} \quad (9)$$

Where

- Within-class matrix

$$S_W = \sum_{i=1}^C S_i \text{ where } S_i = \sum_{x \in C_i} (x^{(i)} - u_i)(x^{(i)} - u_i)^T \text{ and } u_i = \frac{1}{N_i} \sum_{x \in C_i} x^{(i)} \quad (10)$$

- Between-class matrix

$$S_B = \sum_{i=1}^C N_i (u_i - u)(u_i - u)^T \text{ where } u = \frac{1}{n} \sum_{i=1}^n x^{(i)} \quad (11)$$

Thus, if  $S_W$  is a non-singular matrix, and can be inverted, then the Fisher's criterion is maximized when the projection matrix  $\Theta^*$  is composed of the eigen-vectors of:

$$S_W^{-1} S_B \quad (12)$$

Noticed that, there will be at most  $C-1$  eigenvectors with non-zero real corresponding eigenvalues  $\lambda_i$ . This is because  $S_B$  is of rank  $(C-1)$  or less. So we can see that LDA can represent dimensionality reduction of the problem as a classical PCA.

## 2.2 Theory about Kernel Discriminant Analysis

To kernelize the LDA, we consider the problem in a feature space  $\mathcal{F}$  induced by some non-linear mapping function  $\phi : \mathbb{R}^d \rightarrow \mathcal{F}$ . An inner product  $\langle \cdot | \cdot \rangle_{\mathcal{F}}$  can be defined in  $\mathcal{F}$  which makes for so called reproducing the kernel Hilbert Space (RKHS). More precisely  $\langle \phi(x_i) | \phi(x_j) \rangle_{\mathcal{F}} = K(x_i, x_j)$  holds where  $K(\cdot, \cdot)$  is a p.s.d. kernel function. Thus, to find the linear discriminant in  $\mathcal{F}$ , we need to maximize:

$$(\mathcal{P}_K) : \quad \max_{\alpha} J_{\phi}(\alpha) = \max_{\alpha} \frac{\alpha^T S_b^{\phi} \alpha}{\alpha^T S_w^{\phi} \alpha} \quad (13)$$

Where

$$\begin{aligned} S_b^{\phi} &= \sum_c n_c (K_c K_c^T - K K^T) \\ S_w^{\phi} &= K^2 - \sum_c n_c K_c K_c^T \\ K_c &= \frac{1}{n_c} \sum_{i \in c} K_{i,j}, \quad K = \frac{1}{n} \sum_i K_{i,j} \end{aligned} \quad (14)$$

In practice  $S_b^{\phi}$  and  $S_w^{\phi}$  are calculated as

$$S_b^{\phi} = K W K, \quad S_w^{\phi} = K K \quad (15)$$

Where  $K$  is the kernel matrix ( $K_{i,j} = K(x_i, x_j)$ ) and  $W$  is defined as

$$W_{i,j} = \frac{\mathbb{I}(x_i = x_j = k)}{\text{Card}(\mathcal{C}_k)} \quad (16)$$

According to the equation (6), one get the  $\alpha$  which maximizes the objective above, denote by  $\alpha^*$  this solution. Thus,  $\omega_*$  in the initial space is given by:

$$\omega_* = \sum_i \alpha_i^* \phi(x_i) \quad (17)$$

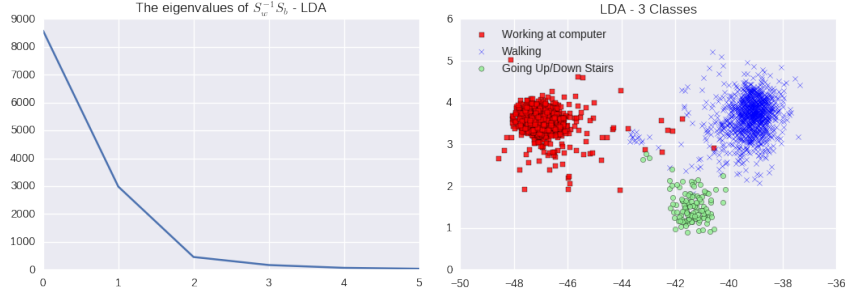
And finally using the kernel trick, we don't need to compute explicitly  $\phi$ , and the projection of a new sample  $x$  over  $\omega_*$  is given by:

$$\omega_*^T \phi(x) = \sum_i \alpha_i^* K(x_i, x) \quad (18)$$

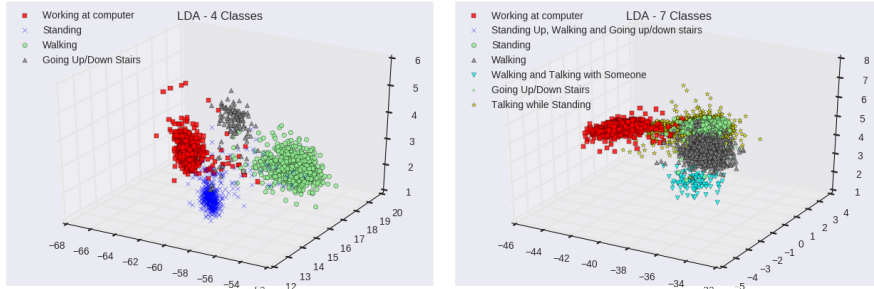
## 2.3 LDA and KDA results

### 2.3.1 LDA results

The figures 4(a) and 4(b) presents the results of applying our implementation of LDA to the data set after the preprocessing described above. We notice that the different classes are separable in the projection space, in particular, for 3 classes (resp. 4 classes), the projected data are separable in a plane (resp. 3D space). Thus after the LDA step, one can apply a classical classifier as Logistic-Regression or SVM to classifier the data in the projection space. The table 2 presents the classification accuracy of the Logistic-Regression and the SVM classifiers after projecting the data with LDA, we notice that we get a better accuracy for LDA+SVM comparing to simply applying SVM to the raw data.



(a) The first eigenvalues of  $S_w^{-1} S_b$  (left) and result of our implementation of LDA applied to 3 classes (right).



(b) Result of our implementation of LDA applied to 4 classes (left) and 7 classes (right). For 7 classes we have plotted only the three first components of the projected data.

### 2.3.2 KDA results

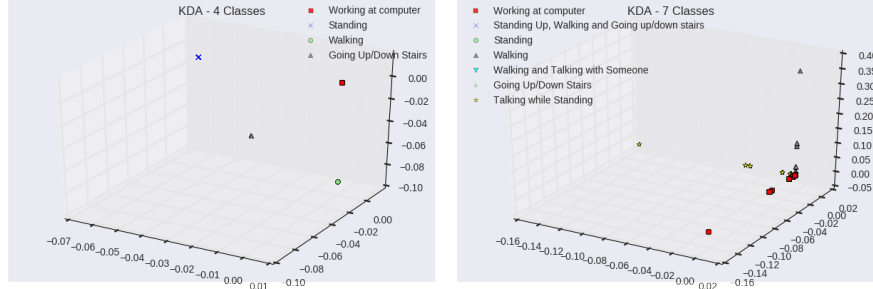
The figures 4(c) and 4(d) presents the results of applying our implementation of KDA to the data set after the preprocessing described above. We notice that the different classes are perfectly separable in the projection space, in particular, for 3 classes (resp. 4 classes), each class is projected in a point of the plane (res. in a point of the 3D space). Thus after the KDA step, as LDA, one can apply a classical classifier as Logistic-Regression or SVM to classifier the data in the projection space. The table 2 presents the classification accuracy of the Logistic-Regression and the SVM classifiers after projecting the data with KDA, we notice that we get a better accuracy for LDA+Log-Reg comparing to applying simply Log-Reg to the raw data. And we get the better accuracy by combining KDA and SVM, ie applying the SVM classifier to the projected data with KDA (SVM+KDA : 99% accuracy).

## 3 Other single frame approaches

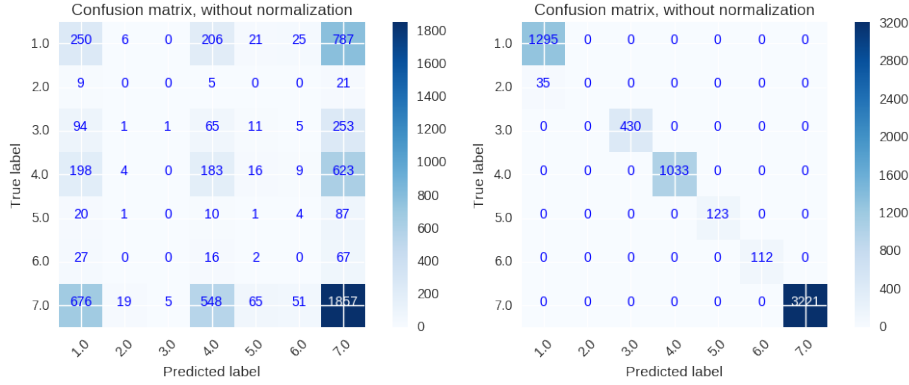
In addition to LDA and KDA, we tried other single frame approaches using popular machine learning algorithms. We used the following algorithms: Logistic regression, support vector machines (SVM),



(c) The first eigenvalues of  $S_w^{-1} S_b$  (left) and result of our implementation of KDA applied to 3 classes (right).



(d) Result of our implementation of KDA applied to 4 classes (left) and 7 classes (right). For 7 classes we have plotted only the three first components of the projected data.



(e) Confusion matrices of LDA (left) and KDA (right).

k-nearest neighbors (knn) and gradient boosting. The hyperparameters of these algorithms were selected by cross validation among a grid of possible values.

To make a fair comparison, we made the same preprocessing on the raw data and thus we used the same features as we did for LDA and KDA. Here are the results we got on the classification accuracy score calculated by cross validation. For Logistic Regression and SVM, we can see that the scores we got slightly lower than the scores obtained on the projections done by LDA and KDA.

Table 1: Classification accuracy score of classical machine learning classifiers.

Algorithm	Classification accuracy score
Logistic Regression	0.884
SVM	0.897
Gradient Boosting	0.913
KNN	0.888



Table 2: Classification accuracy of LDA and KDA with LogReg and SVM.

Algorithm	Classification accuracy score
LDA+Logistic Regression	0.884
LDA+SVM	0.901
KDA+Logistic Regression	0.925
KDA+SVM	0.994

## 4 An Unsupervised Approach for Automatic Activity Recognition based on Hidden Markov Model Regression [5]

In the proposed model, each activity is represented by a regression model and the switching from one activity to another is governed by a hidden Markov chain. The MHMMR parameters are learned in an unsupervised way from unlabelled raw acceleration data acquired during human activities using the Expectation-Maximization (EM) algorithm where no activity labels are needed.

The acceleration data are presented as multidimensional time series presenting various regime changes. In such context, the goal is to provide an automatic partition of the data into different segments (regimes), each segment being considered afterwards as an activity.

In Hidden Markov Model Regression (HMMR), each time series is represented as a sequence of observed univariate variables  $(y_1, y_2, \dots, y_n)$ , where the observation  $y_i$  at time  $t_i$  is assumed to be generated by the following regression model:

$$y_i = \beta_{z_i} + \sigma_{z_i} \epsilon_i ; \epsilon_i \sim \mathcal{N}(0, 1) , (i = 1, \dots, n) \quad (19)$$

where  $z_i \in (1, \dots, K)$  is a hidden discrete-valued variable. In this application case,  $z_i$  represents the hidden class label (activity) of each acceleration data point and  $K$  corresponds to the number of considered activities. The variable  $z_i$  controls the switching from one polynomial regression model associated to one activity, to another of  $K$  models at time  $t_i$ .

The HMMR assumes that the hidden sequence  $z = (z_1, \dots, z_n)$  is a homogeneous Markov chain of first order parameterized by the initial state distribution  $\pi$  and the transition matrix  $A$ . Regarding the multiple regression case, the model can be formulated as follows:

$$y_i^{(1)} = \beta_{z_i}^{(1)T} + \sigma_{z_i}^{(1)} \epsilon_i \quad (20)$$

$$y_i^{(2)} = \beta_{z_i}^{(2)T} + \sigma_{z_i}^{(2)} \epsilon_i \quad (21)$$

$$\vdots \quad \vdots \quad (22)$$

$$y_i^{(d)} = \beta_{z_i}^{(d)T} + \sigma_{z_i}^{(d)} \epsilon_i \quad (23)$$

$$(24)$$

where  $d$  represents the dimension of the time series.

The Multiple HMMR model is therefore fully parameterized by the parameter vector :

$$\theta = (\pi, A, B_1, \dots, B_k, \Sigma_1, \dots, \Sigma_k) \quad (25)$$

where  $B_k$  is the matrix of the multiple regression model parameters associated with the regime (class)  $z_i = k$  and  $\Sigma$  its corresponding covariance matrix.

The parameters are estimated by maximizing the observed data likelihood through the Expectation-Maximization (EM) algorithm. The log-likelihood to be maximized in this case is written as follows:

$$\mathcal{L}(\theta) = \log p(y_1, \dots, y_n; \theta) = \log \sum_z p(z_1; \pi) \prod_{i=2}^n p(z_i | z_{i-1}; A) \prod_{i=1}^n \mathcal{N}(y_i; B_{z_i}^T, t_i, \Sigma_{z_i}) \quad (26)$$

## Conclusion

To conclude, to get good performances, it is extremely important to do good feature engineering in a machine learning task, and especially when we use classical classifiers such as Logistic Regression and SVM. Indeed, we get some good results accuracy for these classifiers after the preprocessing step. However, we have shown that dimension reduction based LDA improves these results because it projects the data in a projection space where the data are separable. Furthermore, the kernel version of LDA improves significantly the accuracy scores (see KDA+SVM : 99.4%), especially, we have noticed that the different classes are projected in 3D points considering 4 classes, which simplifies the classification task. KDA struggles when classes are mixed. For instance, the label 2 contains 3 different activities in the same time, and it was totally missclassified. This issue can be solved by considering a multi-task classification, means structured output prediction. Unfortunately, our implementation of the probabilistic based model (MRHMM) does not work very well, due to some errors caused by unstable log probabilities computation. We will instead provide its code in a separate notebook.

## References

- [1] Andrea Mannini and Angelo Maria Sabatini. Machine learning methods for classifying human physical activity from on-body accelerometers. *Sensors*, 10(2):1154–1175, 2010.
- [2] Andreas Bulling, Ulf Blanke, and Bernt Schiele. A tutorial on human activity recognition using body-worn inertial sensors. *ACM Computing Surveys (CSUR)*, 46(3):33, 2014.
- [3] Pierluigi Casale, Oriol Pujol, and Petia Radeva. Human activity recognition from accelerometer data using a wearable device. In *Iberian Conference on Pattern Recognition and Image Analysis*, pages 289–296. Springer, 2011.
- [4] Adil Mehmood Khan, Y-K Lee, SY Lee, and T-S Kim. Human activity recognition via an accelerometer-enabled-smartphone using kernel discriminant analysis. In *2010 5th International Conference on Future Information Technology*, pages 1–6. IEEE, 2010.
- [5] Ferhat Attal, Samer Mohammed, Mariam Dedabrishvili, Faicel Chamroukhi, Latifa Oukhellou, and Yacine Amirat. Physical human activity recognition using wearable sensors. *Sensors*, 15(12):31314–31338, 2015.