

Subspace Clustering for Action Recognition with Covariance Representations and Temporal Pruning

Giancarlo Paoletti, Jacopo Cavazza, Cigdem Beyan and Alessio Del Bue

Pattern Analysis and Computer Vision, Istituto Italiano di Tecnologia, via Enrico Melen 83, 16152 Genova, Italy

{giancarlo.paoletti, jacopo.cavazza, cigdem.beyan, alessio.delbue}@iit.it

Abstract—This paper tackles the problem of human action recognition, defined as classifying which action is displayed in a trimmed sequence, from skeletal data. Albeit state-of-the-art approaches designed for this application are all supervised, in this paper we pursue a more challenging direction: Solving the problem with *unsupervised learning*. To this end, we propose a novel subspace clustering method, which exploits covariance matrix to enhance the action’s discriminability and a timestamp pruning approach that allow us to better handle the temporal dimension of the data. Through a broad experimental validation, we show that our computational pipeline surpasses existing unsupervised approaches but also can result in favorable performances as compared to supervised methods.

I. INTRODUCTION

Human Action Recognition (HAR) plays a crucial role in computer vision since related to a broad spectrum of artificial intelligence applications (such as video surveillance, human-machine interaction or self-driving cars to name a few [1]). Given a trimmed sequence, in which a single action or activity is assumed to be present, the final goal of HAR is to correctly classifying it. Although significant progresses have been made in the last years, accurate action recognition in videos is still a challenging task because of the complexity of the visual data e.g., due to varying camera viewpoints, occlusions and abrupt changes in lighting conditions. As an all-in-one solution to these problems, skeleton-based HAR is surely the paradigm to embrace, considering also its beneficial characteristics of being privacy-preserving. In skeleton-based HAR, action/activity sequences are represented through the multi-dimensional time series of joints, located at the intersection of skeletal bones, whose position is tracked in time typically through either motion capture systems or depth sensors.

Recently, skeleton-based HAR has undergone to the same paradigm shift which was registered in other fields of pattern recognition: Hand-crafted data encodings fed into engineered classifiers have been replaced by data-driven feature representation with an end-to-end classification pipeline [2]. Yet, both paradigms leverage a fully *supervised* learning approach to accomplish the task. Each sequence is in fact assumed to be (manually) annotated by the action/activity it involves. Other than being time-consuming and prone to human errors, sequence annotations compromise the scalability to the big data regime. As an alternative, unsupervised approaches seem attractive since they offer an advantage regarding computational and methodological burden, as well as providing an interesting application towards more novel real-life scenarios.

In this work, we consider *subspace clustering* to tackle HAR in a fully unsupervised paradigm. Subspace clustering was first introduced in Computer Vision to segment dynamic moving objects [3], [4] and it postulates that high-dimensional data (here, skeletal joints) can be represented as a union of subspaces, each of them having a much lower dimensionality (i.e. low-rank) and simpler geometrical structure. Each subspace usually corresponds to a class (here, to an action or an activity). The key idea in subspace clustering is to learn encodings that are then used to construct an affinity matrix from which the data can be clustered together according to the modelled (dis)-similarities between samples [5]. Although, this is usually achieved through a self-expressive model in which each data point is expressed as a linear combination of the remaining ones, additional constraints, such as sparsity, were also adopted [6].

Despite the fact that subspace clustering has become a powerful technique for problems such as face clustering or digit recognition, its applicability to the problems like skeleton-based HAR was only explored by a limited number of works [7], [8], [9]. This is due to many operative limitations including how to handle the temporal dimensions, the inherent noise present in the skeletal data and the related computational issues.

In this paper, we propose two alternative computational strategies to help and support subspace clustering methods in handling the temporal dimensions of action sequences. On the one hand, we encode the raw skeletal trajectories using a covariance representation, which has been shown to be effective for the solving HAR problems [10]. Additionally, we propose a computational strategy to prune the instantaneous body poses – termed *timestamps* hereafter – whose temporal aggregation produces an action sequence. As the result of temporal pruning, we are able to select the most representative timestamps, which are exploited to compress the original action sequence to a fixed duration. Consequently, this *temporal pruning* can be adopted as a successful pre-processing step to accommodate for the usage of a subspace clustering method for HAR.

Through a comprehensive experimental analysis, we validate the impact on HAR of covariance representations and temporal pruning. Eventually, we also demonstrate their degree of complementary to the extent that the performance of a fully unsupervised recognition pipeline can be enhanced. Surprisingly, the overall performance of the proposed unsupervised

approaches can almost fill the gap with state of the art supervised methods.

Overall, we deem that our experimental findings would help practitioners in re-thinking the way HAR is approached, raising the attention in the desirable shift towards more agile unsupervised learning frameworks.

II. RELATED WORK

In this Section, we have reviewed the action recognition methods relying on covariance representation, various subspace clustering algorithms as well as the state-of-the-art supervised approaches for skeleton-based HAR.

Subspace clustering. Subspace clustering has been a popular computational framework in the machine learning community as well as the computer vision and image processing communities (e.g., image representation and compression [11], image segmentation [12], motion segmentation [13]). It aims at finding subspaces each containing a group of data points and then performing clustering based on these subspaces [5].

There has been a lot of work presenting many different subspace clustering methods. Most of the subspace clustering methods learn an affinity matrix and then apply spectral clustering, e.g., low-rank representation [14], [15]. Self-representation based subspace clustering methods reconstruct a sample from a linear combination of other samples [6], [14], [16], [17] and they have proven their effectiveness for high-dimensional data. Sparse subspace clustering integrates l_1 -norm regularization, which mostly results in improvements in the clustering performances [6]. The temporal Laplacian regularization was proposed in [8] and also adopted in other works e.g., [9] to better model kinematic data for the sake of action detection and segmentation.

Most existing subspace clustering methods rely on hand-crafted representations. Instead, more powerful representations can be learned through deep learning, which effectively cluster data samples from non-linear subspaces [18]. Deep subspace clustering methods apply embedding and clustering jointly, typically with an autoencoder network e.g., in [18], [19]. This results in an optimal embedding subspace for clustering, which is more effective compared to conventional clustering methods. Deep adversarial subspace clustering methods, on the other hand, learn more effective sample representations using deep learning while exploiting adversarial learning to supervise and, thus, progressively improve the performance of subspace clustering. This is done by using a subspace clustering generator and a quality-verifying discriminator which are adversarially learned against each other.

Covariance encoding for HAR. The idea of encoding 3D-skeleton dynamics within a single hand-crafted kernel representation has been proposed often in HAR. For instances, it has been shown that Hankel matrices can efficiently model action dynamics when used in tandem with a Hidden Markov Model [20] or a Riemannian nearest neighbours with class-prototypes [21]. Lie group [22] and associated Lie algebra [23] can be effective in modelling human actions and activities by means of roto-translations. Likewise, generic deforming bodies can

be efficiently modelled over variations of Stiefel manifolds [24]. Surely, within the class of kernel representations, a major role is played by a specific symmetric and positive definite (SPD) operator: Covariance matrices (COV). Originally envisaged for image classification and detection [25], COV is an effective representation for skeleton-based HAR since capable of modelling second-order statistics. It was used in tandem of a variety of classification pipelines, such as a temporal pyramid [26] or max-margin approaches [27], [28]. Formal studies have tried to enhance the capability of such operators in modelling non-linear correlations among the data [29], [30]. Kernel approximation was recently investigated in order to speed up the computational pipeline and ensure scalability towards the big data regime [31].

Even though prior work focused on the effectiveness of covariance representations applied to supervised learning pipelines, we instead demonstrate its capabilities for unsupervised learning.

State-of-the-art supervised approaches for skeleton-based HAR. The current mainstream paradigm in skeleton-based HAR is the possibility of learning a feature representation from the data itself, in tandem with the final action classifier. As one of the seminal works in this direction, a hierarchy of bidirectional recurrent neural networks is used by [32] to represent in a bottom-up fashion all the structural relationships between body parts (torso, legs, arms) in the human skeleton. Long-Short Term Memory (LSTM) models can be proficiently applied to 3D action recognition [33], [34]. Throughout the years, LSTM networks have been modified to better accommodate for the task: for instance, by applying a novel mixed-norm regularization term and dropout [35] or recurring to attention mechanisms [36]. Alternatively, joint trajectories are casted into colored images by producing the so-called distance maps [37], [38], [39]. By means of them, usual convolutional neural networks such as AlexNet, despite originally proposed for image classification, can be adapted to HAR [37], [38]. Surely, the most active and recent direction of research leverages the possibility of encoding the whole human skeleton as a graph, furthermore processing it through a graph-convolutional neural network [40], [41].

All such approaches can fully exploit the benefits of an end-to-end and data-driven training since relying on a fully supervised regime in which the sequences to be classified are annotated. Differently, in this paper we pursue the more challenging direction of adopting an unsupervised strategy, relying on subspace clustering. Similarly to what done by [42] for auto-encoders and [43] for generative adversarial networks, the goal of this paper is to propose new computational architectures and evaluate their effectiveness in comparison with supervised learning paradigms.

III. METHODOLOGY

In this Section, we present our computational pipeline which is based on covariance representations and timestamps pruning. In order to properly ablate on the relative importance

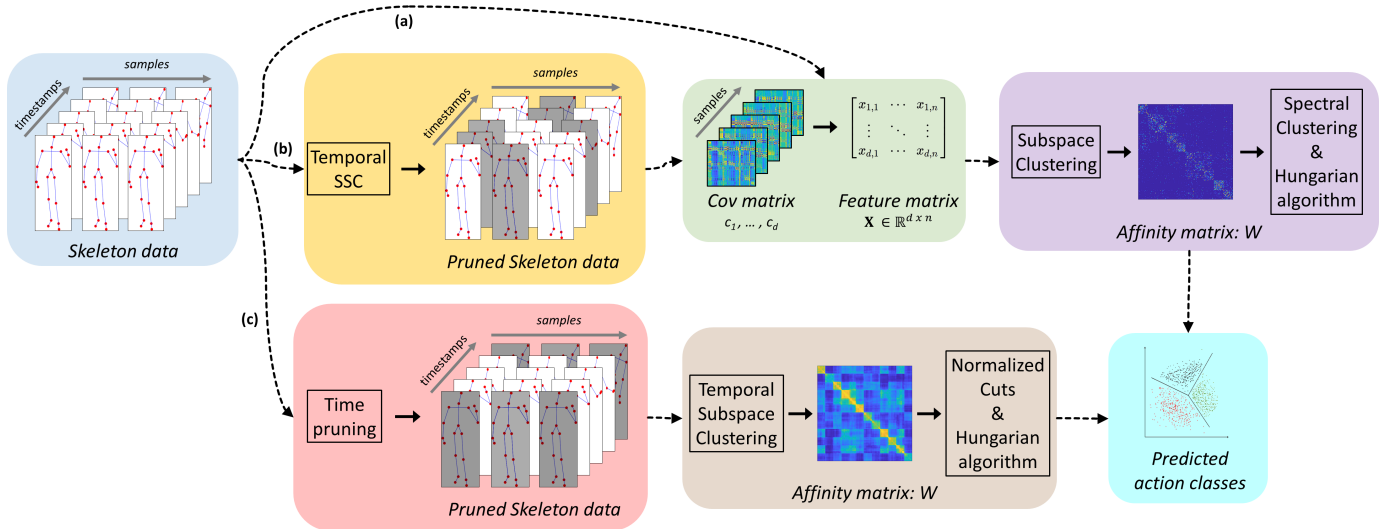


Fig. 1. Pipeline of the proposed unsupervised methods for HAR: (a) A covariance descriptor is applied to each sample. Given the obtained covariance matrix is square and symmetrical, we take only the upper (can be also lower) triangular part including the diagonal and flatten it. This results in a new matrix (X) having size $samples \times features$. Following that, any subspace clustering technique can be applied to obtain an affinity graph matrix W . Then, spectral clustering is applied using W to obtain cluster labels and the Hungarian algorithm finds the matching between the cluster labels (predicted action classes) and the ground-truth labels. (b) The skeletal data of each sample is temporally pruned using temporalSSC and then the pruned data is processed as in (a). (c) Each sample is pruned by using various strategies. Afterwards, temporal subspace clustering is applied to obtain an affinity graph matrix W . The normalized cuts is applied to obtain cluster labels and the Hungarian algorithm matches the cluster labels with the ground-truth labels.

of them, we will consider the following computational variants of the pipeline:

- *Section III-A and Figure 1(a)*. We apply covariance encoding as the descriptor, whose result is given as an input to a subspace clustering method that is based on the self-expressiveness property of data.
- *Section III-B and Figure 1(b)*. We apply the proposed temporal pruning approach (namely *temporalSSC*) as a pre-processing stage while the rest of the pipeline follows the previous setting.
- *Section III-C and Figure 1(c)*. We use the Temporal Subspace Clustering to show the effectiveness of a dictionary-based subspace clustering for temporal series of data when applying temporal regularization on top of the (optional) encoding through covariance.

A. Subspace clustering methods based on self-expressiveness property and covariance representatio

The usage of a covariance representation as the data encoder and the subspace clustering for solving HAR can be described as follows.

Data encoding through covariance representation.

Through either a motion capture system or a depth sensor, an action is represented as the collection in time of K joints 3D positions $\mathbf{p}_1(t), \dots, \mathbf{p}_K(t)$. By using $\mathbf{p}(t)$ to denote the column vectorization of all such 3D positions for a fixed timestamp, we represent an action sequence as the covariance matrix

$$\Lambda = \frac{1}{T} \sum_t (\mathbf{p}(t) - \boldsymbol{\mu})(\mathbf{p}(t) - \boldsymbol{\mu})^\top, \quad (1)$$

where T denotes the number of timestamps and $\boldsymbol{\mu}$ is the temporal average of $\mathbf{p}(t)$.

We then vectorize the covariance matrix through a flattening operation which exploit the property of Λ in being symmetrical. That is, $\Lambda = \Lambda^\top$. Therefore, when flattening, we extract the diagonal elements of Λ (which are Λ_{ii}) and the upper-triangular ones (that is, $\Lambda_{ij}, j > i$). The lower triangular part can be ignored since it is equal to upper triangular one. Such flattening operation casts the $3K \times 3K$ matrix Λ into a $3K \cdot (3K - 1)/2$ column vector. The flattened covariance representation is used as one data point, which then given to the subspace clustering algorithm as the input.

Subspace Clustering. Let us consider a collection of D -dimensional data-points $\mathbf{x}_1, \dots, \mathbf{x}_N$. Subspace clustering [5] attempts to cluster $\mathbf{x}_1, \dots, \mathbf{x}_N$ into groups (termed *subspaces*) which share common geometrical relationships as the well-known *self-expressiveness property*. The problem can be formalised as finding a $N \times N$ matrix C of coefficients such that

$$\mathbf{X} = \mathbf{X}C \text{ subject to } \text{diag}(C) = 0, \quad (2)$$

where \mathbf{X} is the $D \times N$ matrix, which stacks by columns the data points \mathbf{x}_j . The constraint $\text{diag}(C) = 0$ avoids the trivial solution corresponding to C being the identity matrix. Ultimately, the geometrical relationship that we are interested in modelling is a linear relationship in which each data-point can be described as a linear combination. As a consequence of that, the subspaces are linear in turn. The constraint $\text{diag}(C) = 0$ is fundamental to avoid the trivial (and useless) solution $\mathbf{x}_j = \mathbf{x}_j$. Specifically, the self-expressiveness property (2) attempts to estimate each data points as a linear combination of *different data points*. This allows to capture the geometrical inter-dependencies among the data points themselves.

An important aspect regarding subspace clustering is the

way the matrix \mathbf{C} is obtained. A number of works proposed to solve this problem through optimization [44], [6], [45], [46], [47], [18] and different strategies have been adopted to constraint the solution. In subspace segmentation via Least Squares Regression (**SS-LSR**) [44], a Frobenius norm is introduced to promote a L^2 penalty, obtaining

$$\min \|\mathbf{C}\|_F \text{ subject to } \mathbf{X} = \mathbf{X}\mathbf{C}, \text{diag}(\mathbf{C}) = 0. \quad (3)$$

Another popular manner of constraining the coefficient matrix \mathbf{C} is to impose sparsity [6], [47], [18]. As in the Sparse Subspace clustering via Alternating Direction Method of Multipliers (**SSC-ADMM**) [6], the problem formulation is framed as

$$\min \|\mathbf{C}\|_1 \text{ subject to } \mathbf{X} = \mathbf{X}\mathbf{C}, \text{diag}(\mathbf{C}) = 0, \quad (4)$$

while using the alternating direction method of multipliers (ADMM) algorithm to foster convergence by solving a stack of easier sub-problems. As an alternative to ADMM, Sparse Subspace Clustering by Orthogonal Matching Pursuit (**SSC-OMP**) [46] approaches a similar problem with a different optimization technique.

The previous formalism in Eq. (4) was extended in the Deep Subspace Clustering Networks (**DSC-Nets**) [18] by having the hidden layer of an autoencoder implementing either equation (3) or equation (4). The Elastic Net (**EnSC**) [47] approach uses a convex combination of L^2 and L^1 constraint on \mathbf{C} to increase performance, while also boosting the scalability due to the usage of oracle sets to better pre-condition the solution. Dense subspace clustering (**EDSC**) [45] approaches the problem by attempting to apply the self-expressiveness loss on a dictionary which is used to describe the data, while also taking into account outliers.

Once the matrix of coefficient \mathbf{C} is found, an affinity graph matrix \mathbf{W} is built by setting the weights on the edges between the nodes through $\mathbf{W} = \mathbf{C} + \mathbf{C}^\top$. Spectral clustering is later applied to \mathbf{W} to obtain the clustering labels, by assigning each of the N datapoint \mathbf{x}_j into its corresponding subspace. The final step is therefore apply Hungarian algorithm to compare and map subspace labels into actual class labels [5].

B. Temporal pruning via Sparse Subspace Clustering (temporalSSC)

In addition to utilize subspace clustering as an unsupervised learning method to perform action recognition, in this paper, we also exploit such family of techniques to solve another task: temporal pruning. That refers to utilizing subspace clustering on the raw joint coordinates $\mathbf{p}(t)$. Here, different from the previous section, each data point to be clustered is not an action sequence, but a single data point of an action (Figure 1(b)). In other words, rather than applying subspace clustering to group action sequences, we exploit subspace clustering to the group skeletal poses at a given timestamp. Our assumption is that the processed skeleton data might contain similar or even redundant poses over time. To address this, we apply temporal pruning, which potentially captures the similarities over time with respect to the kinematic execution.

A relevant parameter for temporal pruning is the number of subspaces ϕ , which corresponds to the length of the new pruned skeleton data, which was set based on the following strategies.

1) *min ϕ* : the temporal length of the entire dataset is fixed to be equal to the shortest time duration across all the sequences in the skeletal dataset, this is done by using the random permutation of each sample timestamps.

2) *min temporalSSC*: subspace clustering method SSC_ADMM is used to get ϕ equal to the shortest time duration across all the sequences in the skeletal dataset.

3) *percentage temporalSSC*: the temporal length of each sample of the dataset is determined by selecting a percentage value for ϕ (in our experiments we chose to keep the 75%, 50% or 25% of the sample temporal length) and applying temporalSSC.

4) *threshold temporalSSC*: the temporal length of each sample of the dataset is determined by selecting a percentage value for ϕ (in our experiments we chose to keep the 75%, 50% or 25% of the sample temporal length), which is used as a threshold value for temporalSSC. If a certain sample of the dataset has a temporal length superior to ϕ , temporalSSC is therefore applied to match this threshold value.

Once ϕ is fixed according to one of the previous strategy, we can now retrieve all the timestamps t_1, \dots, t_s, \dots assigned to a given subspace. Afterwards, we average the corresponding skeletal positions $\mathbf{p}(t_1), \dots, \mathbf{p}(t_s), \dots$. The so-obtained average skeletal position is adopted to replace the original one and the procedure is iterated across all the different subspaces. For the sake of clarity, let us exemplify the procedure in a particular case. For instance, lets assume that the number of subspaces is set to be $\phi = 2$ and the original action sequence has 5 timestamps to which are associated the following body poses $[\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \mathbf{p}_4, \mathbf{p}_5]$. Once temporalSSC is runned on top of the sequence $[\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \mathbf{p}_4, \mathbf{p}_5]$, let assume that the corresponding output is $[1, 1, 2, 1, 2]$. So, temporalSSC is grouping $\mathbf{p}_1, \mathbf{p}_2$ and \mathbf{p}_4 in a subspaces and $\mathbf{p}_3, \mathbf{p}_5$ in another one. Then, we define the pruned action sequence as $[\mathbf{p}'_1, \mathbf{p}'_2]$, where $\mathbf{p}'_1 = \frac{1}{3}(\mathbf{p}_1 + \mathbf{p}_2 + \mathbf{p}_4)$ and $\mathbf{p}'_2 = \frac{1}{2}(\mathbf{p}_3 + \mathbf{p}_5)$.

Once the temporal pruning is performed, the covariance representation is applied to the new data and subspace clustering is adopted as in Section III-A.

C. Temporal Subspace Clustering based on dictionary and temporal Laplacian Regularization

Even though subspace clustering methods explained in Section III-A build the affinity matrix \mathbf{W} by exploiting the self-expressiveness property of data, they do not explicitly take into account the temporal dimension of time-series data while building the model adopted for HAR. As a solution, temporal regularization was proposed by Temporal Subspace Clustering (**TSC**) [8]. Precisely, given a dictionary $\mathbf{D} \in \mathbb{R}^{d \times r}$ and a coding matrix $\mathbf{Z} \in \mathbb{R}^{r \times n}$, a collection of data points $\mathbf{X} \in \mathbb{R}^{d \times n}$ can be approximately represented as

$$\mathbf{X} \approx \mathbf{D}\mathbf{Z}, \quad (5)$$

where each data point is encoded using a Least Squares regression, and a temporal Laplacian regularization $L(\mathbf{Z})$ function encourages the encoding of the sequential relationships in time-series data. This can be done by minimizing

$$\min_{\mathbf{Z}, \mathbf{D}} \|\mathbf{X} - \mathbf{DZ}\|_F^2 + \lambda_1 \|\mathbf{Z}\|_F^2 + \lambda_2 L(\mathbf{Z}), \quad (6)$$

subject to $\mathbf{Z} \geq 0, \mathbf{D} \geq 0,$

by using the ADMM algorithm to encourage convergence by solving a stack of easier sub-problems. Different from Section III-A, the affinity graph matrix \mathbf{W} is given by the coding matrix \mathbf{Z} by using $\mathbf{W}(i, j) = \frac{z_i^\top z_j}{\|z_i\|_2 \|z_j\|_2}$, since the within-cluster samples (for example the sequential neighbors of a time-series datapoint) are always highly correlated to each other [48], [49]. As final steps of the pipeline, the standard Normalized Cuts [50] and Hungarian algorithms determine the clustering labels necessary for evaluation against the ground-truth.

In Sections III-A and III-B, a (flattened) covariance representation was adopted to encode the actions' kinematics. Computationally, this operation was able to cast an action sequence with a variable temporal duration into a fixed-size embedding which was passed in input to subspace clustering methods based on the self-expressiveness property. Here, differently, TSC leverages a dictionary learning framework which, together with the temporal regularization, should be effective in capturing the temporal variability of the data. To understand to which extent this is true, we would like to *intentionally* get rid of covariance representations within our computational pipeline in order to separately evaluate this two alternative strategies of handling the temporal dimensions of the data.

TSC approach is combined with the following pruning strategies such that a constant temporal length ϕ for all the dataset in use is set as:

1) *TSC min*: the temporal length ϕ of the entire dataset is fixed to be equal to the shortest time duration across all the sequences in the skeletal dataset, this is done by using the random permutation of each timeframe.

2) *TSC max*: the opposite process of *TSC min*. For each instance, its timeframes are replicated until the temporal length ϕ is equal to the longest time duration across all the sequences in the skeletal dataset.

3) *temporalSC + TSC*: spectral clustering is used to get ϕ equal to the shortest time duration across all the sequences in the skeletal dataset.

4) *temporalKm + TSC*: k-means clustering is used to get ϕ equal to the shortest time duration across all the sequences in the skeletal dataset.

IV. 3D ACTION RECOGNITION DATASET

There exists a consistent variability in every HAR dataset due to the length in the performed actions and their complexity, the number of action classes and the technology that was used to capturing them. Prior to experimental analysis, a pre-processing step is performed [20], [21], [22], [23], [28], [30],

[34] in order to fix one root joint located at the hip center, and compute the relative differences of all other $J - 1$ 3D joint positions. This pre-processing is performed at any timestamps $t = 1, \dots, T$ to obtain a $3(J - 1)$ -dimensional (column) vector $p(t)$ of the relative displacements. We used the following dataset for our experimental analysis.

Florence3D (F3D) [51]: a 9-class action dataset (*answer phone, bow, clap, drink, read watch, sit down, stand up, tight lace, wave*) captured using a Microsoft Kinect camera. The actions were performed for two/three times by 10 subjects, resulting in 215 data samples.

UTKinect-Action3D (UTK) [52]: a 10-class action dataset (*carry, clap hands, pick up, pull, push, sit down, stand up, throw, walk, wave hands*) captured using a single stationary Microsoft Kinect camera. Each action was performed for two times by 10 subjects, resulting in 199 data samples.

MSR 3D Action Pairs (MSRP) [53]: includes 12 actions in pairs (*pick up box, put down box, lift box, place box, push chair, pull chair, wear hat, take off hat, put on backpack, take off backpack, stick poster, remove poster*). Each pair has similar features but their relation in terms of motion and shape is different. The actions were performed for three times by 10 subjects, resulting in 353 activity samples.

MSR Action 3D (MSRA) [54]: a 20-class action dataset (*bend, draw circle, draw tick, draw x, forward kick, forward punch, golf swing, hand catch, hand clap, hammer, high arm wave, high throw, horizontal arm wave, jogging, pick up and throw, sideboxing, side kick, tennis serve, tennis swing, two-handwave*) captured by a depth-camera. Each action was performed for three times by 10 subjects, resulting in 557 data samples.

Gaming 3D (G3D) [55]: a 20-class gaming actions dataset (*aim and fire gun, clap, climb, crouch, defend, flap, golf swing, jump, kick left, kick right, punch left, punch right, run, steer a car, tennis swing backhand, tennis swing forehand, tennis serve, throw bowling ball, wave, walk*) captured using a Kinect camera. The actions were repeated for seven times by 10 subjects, resulting in 663 activity samples.

HDM05 [56]: due to class imbalance of the original dataset, we select 14 classes (**HDM-05-14**, *clap above head, deposit floor, elbow to knee, grab high, hop both legs, jog, kick forward, lie down on floor, rotate both arms backward, sit down chair, sneak, squat, stand up, throw basketball*, following the protocol of [27], [30]), and 65 classes (**HDM-05-65**, following the protocol of [57] by grouping together similar actions). The sequences were captured using VICON cameras, resulting in 686 data samples for the former and 2343 data samples for the latter.

MSRC-Kinect12 (MSRC) [58]: a 12-class gesturing dataset, grouped into iconic and metaphoric gestures (*beat both, bow, change weapon, duck, goggles, had enough, kick, lift outstretched arms, push right, shoot, throw, wind it up*). Highly corrupted actions were removed following the protocol as in [26], resulting in 5881 data samples.

TABLE I

CLUSTERING ACCURACY (%) OF SUBSPACE CLUSTERING METHODS AS WELL AS K-MEANS (Km) AND SPECTRAL CLUSTERING (Sc). AVG AND STD STAND FOR THE AVERAGE AND STANDARD DEVIATION OF RESULTS AT EACH COLUMN. THE BEST PERFORMANCE FOR EACH DATASET EMPHASIZED IN BOLD.

Dataset	Km	Sc	EDSC	OMP	DSCN	LSR	SSC	EnSC
F3D	45,58	66,05	54,42	61,40	57,02	60,47	69,12	70,23
UTK	34,67	66,83	52,71	58,79	69,35	57,79	73,97	78,90
MSRP	42,78	52,69	51,90	50,14	49,26	47,31	49,60	49,86
MSRA	41,11	65,17	52,69	43,99	59,91	54,40	57,27	62,84
G3D	31,22	64,71	44,48	45,70	62,59	64,25	65,16	72,25
HDM-05-14	32,36	53,35	52,42	47,67	56,27	51,60	49,13	56,00
HDM-05-65	31,41	44,46	44,43	36,07	30,95	42,98	35,98	42,38
MSRC	61,54	84,34	81,30	51,20	71,35	87,04	62,27	83,27
AVG	40,08	62,22	54,29	49,37	57,09	58,23	57,81	64,46
STD	09,63	11,28	10,82	07,58	11,92	12,66	11,63	13,38

V. EXPERIMENTAL ANALYSIS

In this section we evaluate the methodologies presented in Section III, applied on the dataset reported in Section IV: The Subspace Clustering methods based on self-expressiveness property of data (Section III-A), the Temporal pruning via Sparse Subspace Clustering (temporalSSC) (Section III-B) and the Temporal Subspace Clustering based on dictionary and temporal Laplacian Regularization (Section III-C).

Error metrics and performance evaluation. To monitor the performance in HAR, we will take advantage of classification accuracy defined as

$$ACC(\%) = \left(1 - \frac{\# \text{ of misclassified labels}}{\# \text{ of total labels}}\right) \times 100 \quad (7)$$

and expressed as a percentage. As explained in Section III, the clustering labels are obtained through either spectral clustering [59] or Normalized Cut [50]. Finally, the Hungarian algorithm [5] maps cluster labels into the ground-truth ones.

A. Subspace Clustering methods based on self-expressiveness property of data

In this section we experimentally validate the computational method presented in Section III-A and visualised in Figure 1(a): In order to exploit the self-expressiveness property of data and to encode their temporal information, we implement the covariance descriptor to encode the raw data.

Following that, we used the state-of-the-art subspace clustering methods that are based on the self-expressiveness property of the data to obtain the affinity matrix \mathbf{W} . These methods are: EDSC [45], OMP [46], DSCN [18], LSR [44], SSC [6], EnSC [47] which are described in Section III-A.

Once the coefficient matrix \mathbf{C} and the affinity graph matrix \mathbf{W} were found, spectral clustering and Hungarian algorithm were applied to map the subspace label with the actual class labels [5] as illustrated in Figure 1(a).

Additionally, as a baseline method, we considered two of the most popular clustering method: K-means clustering (**Km**) and spectral clustering (**Sc**) [59] and all the corresponding results are reported in Table I. The overall best performing method is Elastic net Subspace Clustering (EnSC) [47], which ranked highest for five of the nine datasets. For three of these

TABLE II

CLUSTERING ACCURACY (%) OF TEMPORALSSC COMBINED WITH DIFFERENT STRATEGIES AND WHEN STANDARD SSC APPLIED FOR THE FINAL CLUSTERING. THE FIRST COLUMN SHOWS THE SSC'S PERFORMANCES ALONE. AVG AND STD STAND FOR THE AVERAGE AND STANDARD DEVIATION OF RESULTS AT EACH COLUMN. BEST PERFORMANCE OF EACH DATASET EMPHASIZED IN BOLD.

Dataset	SSC	min	min	percentage	ϕ	threshold	ϕ
		ϕ	temporalSSC	temporalSSC		temporalSSC	
F3D	69,12	67,91	66,51	65,12	75%	68,84	50%
UTK	73,97	64,82	80,90	68,34	25%	72,86	75%
MSRP	49,60	48,88	47,88	50,42	25%	49,58	25%
MSRA	57,27	59,61	57,09	62,66	25%	63,02	75%
G3D	65,16	64,86	64,10	69,68	75%	71,49	75%
HDM-05-14	49,13	63,12	59,04	59,33	25%	59,77	25%
HDM-05-65	35,98	41,31	44,00	43,66	25%	41,53	50%
MSRC	62,27	83,79	83,62	83,41	75%	83,14	75%
AVG	57,81	61,79	62,89	62,83		63,78	
STD	11,63	11,90	13,23	11,40		12,53	

five, i.e., UTKinect, MSRAAction3D, and G3D datasets, EnSC's performance is approximately 5% better than the second best performing method.

B. Temporal pruning via Sparse Subspace Clustering (temporalSSC)

This pipeline is similar to (A) but we applied to raw data, before the encoding of the covariance descriptor, different pruning strategies for the temporal dimension of data by using SSC (see Figure 1(b)). For the subspace clustering implementation, we decided to use SSC for its computational efficiency and rapid convergence time.

Table II reports the clustering accuracy of different temporalSSC strategies, along with SSC results of Table I [6] as a baseline comparison. Results of *percentage temporalSSC* and *threshold temporalSSC* are related to the best accuracy along the different percentage values of ϕ (i.e. 75%, 50% and 25%). Only with the exception of F3D (due to its original low dimensionality of the dataset and the extreme pruning of timestamps), the results show that applying temporalSSC overall contributes positively to the clustering performance of SSC [6]: The performance improvement is up to an average 8% among all dataset, where on MSRC (the biggest dataset available) the improvement goes up to 21%.

C. Temporal Subspace Clustering based on dictionary and temporal Laplacian Regularization

Table III reports the unsupervised clustering accuracy of the approach given in Section III-C (as well as illustrated in Figure 1(c)) is applied. We also *TSCmin*, *TSCmax*, *temporalSC* + *TSC*, *temporalKm* + *TSC* with and without covariance descriptor. The last column of that table reports the state-of-the-art performance obtained for each dataset. It is important to highlight that the corresponding state-of-the-art methods are all supervised while all other results given in that table are unsupervised.

The results show that the application of TSC gives the best overall accuracy among all techniques adopted in this paper, Table III demonstrates that the average of results of each implementation (column) is over over 85% among all cases. Except G3D and HDM-05-65 datasets, the average

TABLE III
CLUSTERING ACCURACY (%) OF TSC COMBINED WITH DIFFERENT STRATEGIES OF UNIFORMING TEMPORAL DIMENSION OF EACH DATASET. THE SUPERVISED STATE-OF-THE-ART (S.O.T.A) RESULTS ARE ALSO GIVEN. AVG AND STD STAND FOR THE AVERAGE AND STANDARD DEVIATION OF RESULTS AT EACH COLUMN. BEST UNSUPERVISED PERFORMANCE OF EACH DATASET EMPHASIZED IN BOLD.

Dataset	TSCmin	cov TSCmin	TSCmax	cov TSCmax	temporalSC + TSC	temporalSC + TSC cov	temporalKm + TSC	temporalKm + TSC cov	supervised s.o.t.a.
F3D	84,65	81,40	94,88	81,86	95,81	88,84	87,91	87,44	99,07 [60]
UTK	93,97	96,98	99,50	92,96	96,98	96,98	93,47	83,92	100,00 [61]
MSRP	93,48	81,30	98,02	84,70	88,67	76,20	96,32	71,10	95,50 [10]
MSRA	87,18	79,89	85,64	83,30	82,47	81,13	88,51	87,61	97,40 [10]
G3D	88,99	90,20	85,07	92,61	90,20	92,46	88,84	92,91	96,02 [62]
HDM-05-14	89,80	86,73	80,32	83,82	88,48	84,84	83,97	81,63	99,10 [10]
HDM-05-65	70,51	83,57	75,97	85,62	72,13	84,64	68,42	86,00	96,92 [63]
MSRC	97,96	91,09	99,08	99,05	98,81	97,42	99,00	91,07	98,50 [10]
AVG	88,32	86,40	89,81	87,99	89,19	87,81	88,31	85,21	
STD	7,79	5,59	8,62	5,72	8,18	7,05	8,80	6,31	

accuracy of each method without covariance (*cov*) descriptor is approximately 2% better than a method with *cov* descriptor. The comparisons between the temporal frames selection approaches show that in 5-out-of-8 datasets the pruning of data, therefore reduction of its temporal dimension, is beneficial to encode and represent this type datasets. Whereas, for the datasets MSRP and MSRC, augmenting the data in temporal dimension leads to performance levels better than the state-of-the-art methods, which are all supervised.

VI. CONCLUSION

Human Activity Recognition (HAR) is a challenging problem, which has been solved with different methodologies and the sharp majority of them apply a supervised learning paradigm. This paper particularly focuses on skeletal data analysis and, differently, embraces a fully unsupervised approach to tackle HAR. In this study, we propose a novel clustering pipeline, which combines covariance descriptors and subspace clustering applied to 1) temporally prune the input data and 2) group together similar activities based on their respective category. The aim of temporal pruning is to discriminate better the action sequences that are recognized with an unsupervised method.

The experimental analysis is validated on eight different dataset, which are different from each other in terms of action types, the number of action classes involved as well as the experimental protocol they were captured. Across such a wide variety of experimental benchmarks, our findings show that our proposed pipeline is superior to previous subspace clustering methods relying on the self-expressiveness property of data.

Subspace clustering methods based on the self-expressiveness property can remarkably enhanced in performance by covariance representation to the point that other baseline methods are systematically outperformed. On the other hand, temporal subspace clustering method that relies on dictionary learning and temporal Laplacian regularization combined within our pipeline results in remarkably good HAR performances: This demonstrates the benefits of pruning action sequences along the temporal dimension. Overall, the combination of our experimental findings enable a fully unsupervised pipeline for HAR to

always reduce the gap with supervised approaches, while surprisingly outperforming them in some cases.

REFERENCES

- [1] R. Poppe, "A survey on vision-based human action recognition," *Image and Vision Computing*, vol. 28, no. 6, pp. 976 – 990, 2010.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105.
- [3] C. W. Gear, "Multibody grouping from motion images," *IJCV*, vol. 29, no. 2, pp. 133–150, 1998.
- [4] J. P. Costeira and T. Kanade, "A multibody factorization method for independently moving objects," *IJCV*, vol. 29, no. 3, pp. 159–179, 1998.
- [5] R. Vidal, "Subspace clustering," *IEEE Signal Processing Magazine*, vol. 28, no. 2, pp. 52–68, 2011.
- [6] E. Elhamifar and R. Vidal, "Sparse subspace clustering: Algorithm, theory, and applications," *IEEE TPAMI*, vol. 35, no. 11, pp. 2765–2781, 2013.
- [7] H. Zhang and O. Yoshie, "Improving human activity recognition using subspace clustering," in *International Conference on Machine Learning and Cybernetics*, vol. 3. IEEE, 2012, pp. 1058–1063.
- [8] S. Li, K. Li, and Y. Fu, "Temporal subspace clustering for human motion segmentation," in *Proceedings of IEEE ICCV*, 2015, pp. 4453–4461.
- [9] L. Clopton, E. Mavroudi, M. Tsakiris, H. Ali, and R. Vidal, "Temporal subspace clustering for unsupervised action segmentation," *CSMR REU*, pp. 1–7, 2017.
- [10] J. Cavazza, P. Morerio, and V. Murino, "Scalable and compact 3d action recognition with approximated rbf kernel machines," *Pattern Recognition*, vol. 93, pp. 25–35, 2019.
- [11] W. Hong, J. Wright, K. Huang, and Y. Ma, "Multiscale hybrid linear models for lossy image representation," *IEEE TIP*, vol. 15, no. 12, pp. 3655–3671, Dec 2006.
- [12] A. Y. Yang, J. Wright, Y. Ma, and S. S. Sastry, "Unsupervised segmentation of natural images via lossy data compression," *CVIU*, vol. 110, no. 2, pp. 212 – 225, 2008.
- [13] X. Fan and R. Vidal, "The space of multibody fundamental matrices: Rank, geometry and projection," in *Proceedings of International Conference on Dynamical Vision*. Berlin, Heidelberg: Springer-Verlag, 2006, p. 117.
- [14] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust recovery of subspace structures by low-rank representation," *IEEE TPAMI*, vol. 35, no. 1, pp. 171–184, Jan 2013.
- [15] C. Lu, J. Feng, Z. Lin, T. Mei, and S. Yan, "Subspace clustering by block diagonal representation," *IEEE TPAMI*, vol. 41, no. 2, p. 487501, Feb. 2019.
- [16] H. Hu, Z. Lin, J. Feng, and J. Zhou, "Smooth representation clustering," in *Proceedings of IEEE CVPR*, June 2014, pp. 3834–3841.
- [17] C.-Y. Lu, H. Min, Z.-Q. Zhao, L. Zhu, D.-S. Huang, and S. Yan, "Robust and efficient subspace segmentation via least squares regression," in *Proceedings of ECCV - Volume Part VII*, ser. ECCV12, 2012, p. 347360.

- [18] P. Ji, T. Zhang, H. Li, M. Salzmann, and I. Reid, "Deep subspace clustering networks," in *Advances in Neural Information Processing Systems*, 2017, pp. 24–33.
- [19] X. Yang, C. Deng, F. Zheng, J. Yan, and W. Liu, "Deep spectral clustering using dual autoencoder network," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019, pp. 4061–4070.
- [20] L. Lo Presti, M. La Cascia, S. Sclaroff, and O. Camps, "Gesture modeling by Hanklet-based hidden Markov model," in *Asian Conference on Computer Vision (ACCV)*, 2014.
- [21] X. Zhang, Y. Wang, M. Gou, M. Szaier, and O. Camps, "Efficient temporal sequence comparison and classification using gram matrix embeddings on a riemannian manifold," in *Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [22] R. Vemulapalli, F. Arrate, and R. Chellappa, "Human action recognition by representing 3d skeletons as points in a lie group," in *Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [23] R. Vemulapalli and R. Chellappa, "Rolling rotations for recognizing human actions from 3d skeletal data," in *Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [24] A. Del Bue, J. Xavier, L. Agapito, and M. Paladini, "Bilinear modeling via augmented lagrange multipliers (balm)," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 8, pp. 1496–1508, 2011.
- [25] O. Tuzel, F. Porikli, and P. Meer, "Region covariance: A fast descriptor for detection and classification," in *European Conference on Computer Vision (ECCV)*, 2006.
- [26] M. Hussein, M. Torki, M. Gowayed, and M. El-Saban, "Human action recognition using a temporal hierarchy of covariance descriptors on 3d joint locations," in *International Joint Conference on Artificial Intelligence*, 2013.
- [27] L. Wang, J. Zhang, L. Zhou, C. Tang, and W. Li, "Beyond covariance: Feature representation with nonlinear kernel matrices," in *International Conference on Computer Vision (ICCV)*, 2015.
- [28] P. Koniusz, A. Cherian, and F. Porikli, "Tensor representation via kernel linearization for action recognition from 3d skeletons," in *European Conference on Computer Vision (ECCV)*, 2016.
- [29] M. Harandi, M. Salzmann, and F. Porikli, "Bregman divergences for infinite dimensional covariance matrices," in *Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [30] J. Cavazza, A. Zunino, M. San Biagio, and V. Murino, "Kernelized covariance for action recognition," in *International Conference on Pattern Recognition (ICPR)*, 2016.
- [31] J. Cavazza, P. Morerio, and V. Murino, "Scalable and compact 3d action recognition with approximated rbf kernel machines," *Pattern Recognition*, vol. 93, pp. 25–35, 2019.
- [32] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [33] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "NTU RGB+D: A large scale dataset for 3D human activity analysis," in *Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [34] J. Liu, A. Shahroudy, D. Xu, and G. Wang, "Spatio-temporal LSTM with trust gates for 3D human action recognition," in *European Conference on Computer Vision (ECCV)*, 2016.
- [35] W. Zhu, C. Lan, J. Xing, W. Zeng, Y. Li, L. Shen, and X. Xie, "Co-occurrence feature learning for skeleton based action recognition using regularized deep LSTM networks," in *AAAI Conference on Artificial Intelligence*, 2016.
- [36] J. Liu, G. Wang, P. Hu, L.-Y. Duan, and A. C. Kot, "Global context-aware attention lstm networks for 3d action recognition," in *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [37] P. Wang, Z. Li, Y. Hou, and W. Li, "Action recognition based on joint trajectory maps using convolutional neural networks," in *ACM Multimedia*, 2016.
- [38] C. Li, Y. Hou, P. Wang, and W. Li, "Joint distance maps based action recognition with convolutional neural network," in *IEEE Signal Processing Letters*, 2017.
- [39] Q. Ke, M. Bennamoun, S. An, F. Sohel, and F. Boussaid, "A new representation of skeleton sequences for 3d action recognition," in *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [40] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Two-stream adaptive graph convolutional networks for skeleton-based action recognition," in *Proceedings of IEEE CVPR*, 2019, pp. 12 026–12 035.
- [41] C. Wu, X.-J. Wu, and J. Kittler, "Spatial residual layer and dense connection block enhanced spatial temporal graph convolutional network for skeleton-based action recognition," in *Proceedings of IEEE ICCVW*, 2019, pp. 0–0.
- [42] K. Su, X. Liu, and E. Shlizerman, "Predict & cluster: Unsupervised skeleton based action recognition," *arXiv preprint arXiv:1911.12409*, 2019.
- [43] N. Zheng, J. Wen, R. Liu, L. Long, J. Dai, and Z. Gong, "Unsupervised representation learning with long-term dynamics for skeleton based action recognition," in *AAAI Conference on Artificial Intelligence*, 2018.
- [44] C.-Y. Lu, H. Min, Z.-Q. Zhao, L. Zhu, D.-S. Huang, and S. Yan, "Robust and efficient subspace segmentation via least squares regression," in *ECCV*. Springer, 2012, pp. 347–360.
- [45] P. Ji, M. Salzmann, and H. Li, "Efficient dense subspace clustering," in *Proceedings of IEEE WACV*. IEEE, 2014, pp. 461–468.
- [46] C. You, D. Robinson, and R. Vidal, "Scalable sparse subspace clustering by orthogonal matching pursuit," in *Proceedings of IEEE CVPR*, 2016, pp. 3918–3927.
- [47] C. You, C.-G. Li, D. P. Robinson, and R. Vidal, "Oracle based active set algorithm for scalable elastic net subspace clustering," in *Proceedings of IEEE CVPR*, 2016, pp. 3928–3937.
- [48] S. Li and Y. Fu, "Low-rank coding with b-matching constraint for semi-supervised classification," in *International Joint Conference on Artificial Intelligence*, 2013.
- [49] —, "Learning balanced and unbalanced graphs via low-rank coding," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 5, pp. 1274–1287, 2014.
- [50] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE TPAMI*, vol. 22, no. 8, pp. 888–905, 2000.
- [51] L. Seidenari, V. Varano, S. Berretti, A. Bimbo, and P. Pala, "Recognizing actions from depth cameras as weakly aligned multi-part bag-of-poses," in *Proceedings of IEEE CVPRW*, 2013, pp. 479–485.
- [52] L. Xia, C. Chen, and J. Aggarwal, "View invariant human action recognition using histograms of 3d joints," in *Proceedings of IEEE CVPRW*. IEEE, 2012, pp. 20–27.
- [53] O. Oreifej and Z. Liu, "Hon4d: Histogram of oriented 4d normals for activity recognition from depth sequences," in *Proceedings of IEEE CVPR*, 2013, pp. 716–723.
- [54] W. Li, Z. Zhang, and Z. Liu, "Action recognition based on a bag of 3d points," in *Proceedings of IEEE CVPRW*. IEEE, 2010, pp. 9–14.
- [55] V. Bloom, V. Argyriou, and D. Makris, "Hierarchical transfer learning for online recognition of compound actions," *CVIU*, vol. 144, pp. 62–72, 2016.
- [56] M. Müller, T. Röder, M. Clausen, B. Eberhardt, B. Krüger, and A. Weber, "Documentation mocap database hdm05," Universität Bonn, Tech. Rep. CG-2007-2, June 2007.
- [57] K. Cho and X. Chen, "Classifying and visualizing motion capture sequences using deep neural networks," *CoRR*, vol. 1306.3874, 2014.
- [58] S. Fothergill, H. Mentis, P. Kohli, and S. Nowozin, "Instructing people for training gestural interactive systems," in *SIGCHI Conference*, 2012, pp. 1737–1746.
- [59] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, "Scikit-learn: Machine learning in python," *Journal of machine learning research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [60] C. Li, Z. Cui, W. Zheng, C. Xu, and J. Yang, "Spatio-temporal graph convolution for skeleton based action recognition," in *AAAI Conference on Artificial Intelligence*, 2018.
- [61] X. Zhang, Y. Wang, M. Gou, M. Szaier, and O. Camps, "Efficient temporal sequence comparison and classification using gram matrix embeddings on a riemannian manifold," in *Proceedings of IEEE CVPR*, 2016, pp. 4498–4507.
- [62] P. Wang, W. Li, C. Li, and Y. Hou, "Action recognition based on joint trajectory maps with convolutional neural networks," *Knowledge-Based Systems*, vol. 158, pp. 43–53, 2018.
- [63] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *Proceedings of IEEE CVPR*, 2015, pp. 1110–1118.