

Techniques for automated analysis of saliency in human full-body movement

Serafino Gabriele
Lo Luca

Advisors:
Volpe Gualtiero
Camurri Antonio

2024

Abstract

In contemporary research, the topic of salience has become a focal point, with a growing number of studies dedicated to its exploration. Various fields, including Medical, Financial, Military, Psychology, Kinesiology, have approached the problem in diverse ways.

Our thesis aims to develop effective heuristics and algorithmic approaches for detecting salience in the context of movement. Specifically, we focused on the domain of dance movements, defining a concept of salience within this context.

In our purpose the salience is meant as a change in the behaviour of the dancer: static movement with respect to dynamic movement, repetitive movement with respect to a chaotic one, using an arm and then stopping it for using the other one, and so on....

The study involves a thorough analysis of Motion Capture (MOCAP) data, where we navigate dataset constraints and boundaries to select significant features that aptly represent our chosen domain.

To address the salience detection problem, we employ several approaches: among them, a Supervised Machine Learning (ML) Algorithm.

For the ML approach we defined a specific prototype of Sliding Windows (SW) able to convey information through several frames of the videos, with the goal of giving to the specific model the most significant piece of information in order to have the best results.

Then we tested the algorithm at different levels of generalization.

The goal of this work is to develop a system able to recognize the salience in human full-body context in real time.

Contents

1	Introduction	8
1.1	The problem of automatic salience	8
1.2	Motivation and goals	9
1.3	Thesis structure	9
2	Methodologies	10
2.1	CPD algorithms	10
2.2	Related Work	10
2.2.1	State of the art	13
2.3	Structure of workflow	16
3	Salience definition	19
3.1	Ground Truth	20
4	Low level features	23
4.1	Distal Parts and Head movement - $S_8; S_6, S_7$	24
4.2	Point Density for crouched position - S_4, S_5	25
4.3	Repetitiveness - S_2, S_3	25
4.4	Stage dancer position - S_1	26
4.5	Feature Vector	26
5	Supervised ML approach	28
5.1	Data	31
5.1.1	Data Cleaning	34
5.1.2	Sliding Windows	36
5.1.3	Mid level features	36
5.1.4	Normalization	36
5.2	Samples Creation	36
5.3	Model Selection	36
5.4	Results	36
5.4.1	LOSO	36
5.4.2	LOTO	36
5.4.3	LODO	36
6	Statistic approach	37
6.1	Sliding windows	37
6.2	Mid Level Features	37

6.3	model	37
6.4	Results	37
7	Conclusions	37
7.1	Future Researches	37

List of Figures

1	Cora Gasparotti	10
2	Marianne Gubri	10
3	Muriel Romero	10
4	Camera Qualisys	11
5	Markers Qualisys	12
6	Annotation of the different salience with ELAN.	15
7	One of the patches implemented on EyesWeb for extracting the features.	15
8	Workflow chart	18
9	The red line is over the frame were the salience is annotated (S_4).	21
10	After ≈ 50 frames (1 second) with respect to the salience in Figure 9 the dancer is extended.	22
11	EyesWeb patch for extracting the kinetic energy from a marker	24
12	EyesWeb patch for extracting the Point Density	25
13	EyesWeb patch for extracting the Global kinetic energy.	26
14	Saving the variables which store the features values over the frames on a text file (EyesWeb).	27

List of Tables

Acronyms

3D Three dimensional. 12, 23, 24

CPD Change Point Detection. 16

EST Event Segmentation Theory. 13

LOTO Leave One take Out. 31

ML Machine Learning. 2, 17, 28–30

MOCAP Motion Capture. 2, 11, 12, 16, 31

NaN Not a Number. 34, 35

QTM Qualisys Track Manager. 12

SW Sliding Windows. 2

TSV Tab-Separated Values. 12, 14, 24, 31

1 Introduction

The thesis focuses on defining movement saliency and developing a model for its automatic detection and analysis. The work uses the datasets from the archives of Casa Paganini - Infomus, an international Research Center, and consists of videos of dancers performing on-site captured using a motion capture techniques.

1.1 The problem of automatic salience

Salience comes from the Latin salire, meaning "to leap". Something with salience leaps out at you because it is unique or special in some way.[9][6]

One of the major challenge in automate the analysis of salience in movement is the variety and complexity of human behaviours. First of all the human movements is the result of the combination of different part of the body from the evident wave of a human hand to the slightest grin. Humans are capable of different and diversified movements with a big range of possibilities and this makes salience in movement a subjective interpretation. Human can hide their thoughts overriding our reaction but subtle behaviour is hard to hide it. If we want to analyze this kind of saliency, we are more interested in subtle movements. Instead of we want to analyze a martial artists or a tennis player, we are interested in a wider, bigger range of movements and recognize patterns. And ignore subtle movements that could be just adjusting, or else. So salience in movement is subjective and so really hard to generalize and this makes really hard to develop specific algorithm able to generalize every salience of movement. (salienza in diversi contesti dsono diversi, da aggiungere sopra).

Other than the subjective and diverse context of salience in movement, we need to consider the nature of multidimensional data of movements. Also the data could be subject to noise and variability leading to a hard time to distinguish important pattern to noise.

So to automate the analysis of salience in movement, it's a challenge where it's required a interdisciplinary approach and integrate different techniques. (Da sistemare)

1.2 Motivation and goals

The main reason to undertake this project is to challenge ourselves and try to overcome

la ragione principale per cui abbiamo scelto questa tesi è stata quella di voler sfidare noi stessi in un problema così complesso. Perchè, come ho già definito prima, la risoluzione di questo problema può portare allo sviluppo di tanti metodologie future per la risoluzione di tanti problemi. E volevamo contribuire, anche con un solo step a questo percorso.

Il nostro obiettivo è quello di creare un sistema automatico che analizzi i movimenti umani, tale da individuare quei movimenti definiti salienti. Di conseguenza permettere l'installazione di tale sistema automatico in vari sistemi nel mondo reale dove le macchine potranno catturare e analizzare i movimenti umani senza il bisogno di una supervisione umana.

Understanding the salience in movement and being able to automatically analyze it both in a offline and online way, could be used in many different field. In medical way, if we can detect automatically human movements, we can also detect movement that usually is symptoms of illness. If we can detect movement in sports, the athletes can understand if he is performing useless movements. If we can use it on cameras, the camera can detect if, in the crowd, someone is acting differently or in a harmful way or if its feeling bad. The possibilities of automatic analysis of saliency, can be used to be attach to robots, surveillance system and other vision tools to provide to the user more useful information.

1.3 Thesis structure

2 Methodologies

2.1 CPD algorithms

2.2 Related Work

The material we started from was provided by researchers at Casa Paganini - InfoMus [7]. By exploiting the existing results, our goal was to find algorithms and heuristics for dealing with the problem of salience. The given scenario is the following: three dancers from Casa Paganini - InfoMus [7], Cora Gasparotti (Figure 1), Marianne Gubri (Figure 2) and Muriel Romero (Figure 3), performed dancing movements by emphasizing some specific tasks, in order to highlight a combination of graceful and fluid movements, as well as angular and hurried gestures.

Those performances has been recorded with two professional cameras (frontal and lateral view, 1280×720 , 50fps) at Casa Paganini, and the result of this procedure was a set of *takes* (each *take* has a duration between 30 seconds to 180 seconds).

So, for each dancer, we have between 5 to 15 *takes*.

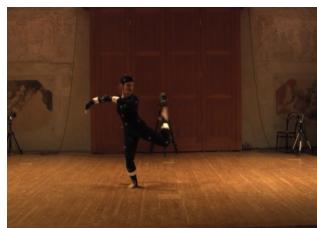


Figure 1:
Cora Gasparotti



Figure 2:
Marianne Gubri



Figure 3:
Muriel Romero

Raw data has been collected by using a Qualisys [8] Motion Capture (MO-CAP) system:

- Qualisys Cameras: these are specialized for capturing and tracking movements in three-dimensional space.

These cameras operate based on the principle of detecting and analyzing markers placed on objects or individuals within their field of view. The cameras emit light, which is then reflected by small, passive markers, such as reflective spheres or retro-reflective devices. This process allows the cameras to discern and meticulously track the position and movement of each marker within their field of view.

In our context, a Qualisys system endowed of 16-cameras was used ($f_s = 100\text{Hz}$). (Figure 4)



Figure 4: Camera Qualisys

- Markers: are integral components of motion capture systems, they are identifiable points placed on objects or the human body. These markers can take various forms, such as reflective spheres, strips, or discs, each designed for specific tracking purposes.
In our cases we used 64 spherical markers positioned on the whole body.
(Figure 5)

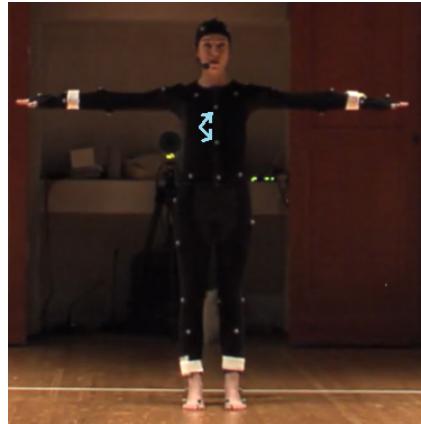


Figure 5: Markers Qualisys

The MOCAP data was collected in a Tab-Separated Values (TSV) file by using Qualisys Track Manager (QTM) [8] which is a software used for files handling MOCAP data, hence for each frame the position of each marker in the 3D space (x,y,z) is available.

2.2.1 State of the art

One of the most important work related to our context was the one done by E. Ceccaldi in her PhD Thesis [2], in that discussion, the problem of salience detection, has been defined with some psychological approach related to the body behaviour, for the definition of the salience, *takes* have been segmented (in the time).

In this work, segmentation is based on boundaries between events, as described by a cognitive theory on how this process unfolds in the human mind, namely the Event Segmentation theory [11]. According to this theory, a boundary is perceived whenever a meaningful (i.e. salient) change is perceived in the ongoing situation. Following the theory, in [2] changes were operationalized as follows:

Thus, salience is defined to a higher level as follows:

- C1. Time: which is the timing and the rhythm of the movement, for instance if the dancer accelerates her actions.
- C2. Space: if the attention and the direction of the movement starts to pointing something different.
- C4. Character Location: if the character moves on the stage, or, in our cases, if the dancer moves from a point to another.
- C6. Causes: Which are the causes and appraisal, if a new state of affairs leads to a subsequent event.

Hence the ground truth was taken by considering these 4 different aspects by psychologists which are trained on EST [11] principles by using a software for segmentation: ELAN [1] (Figure 6).

This software allows to make the segmentation of videos (or an audio track) by defining some classes of event segment and a hierarchy among them.

The *features* which have been considered in that case are the following:

- For the *Time* it was used the General Quantity of Movement and the Chest Quantity of Movement.
- For the *Space* the author employed the Directness of Head Movement.
- The *Character Location* has been detected by looking the Density of Chest Trajectory.

- The causes, instead, has been used only has an additional ground truth.

All these features has been extracted by using an application (patch) developed for EyesWeb XMI (Figure 7), a software platform able to build over the TSV, which contains our raw data, some useful characteristics such as: kinetic Energy, Point density, Directness, and so on...

For instance, kinetic energy (of all the markers) was taken as a measure of the overall amount of movement.

The Chest Quantity of Movement was evaluated by considering the kinetic energy of the marker related to the chest. The Directness of Head is a measurement of how the dancer's head moves in curvilinear trajectory, the higher is the value, the more the movement following a straight line.

For finishing, the Density of chest trajectory is an indicator of whether movement is localized in a small region in the space rather than spanning the whole space, i.e., higher density indicates that the actor has moved in a smaller region.

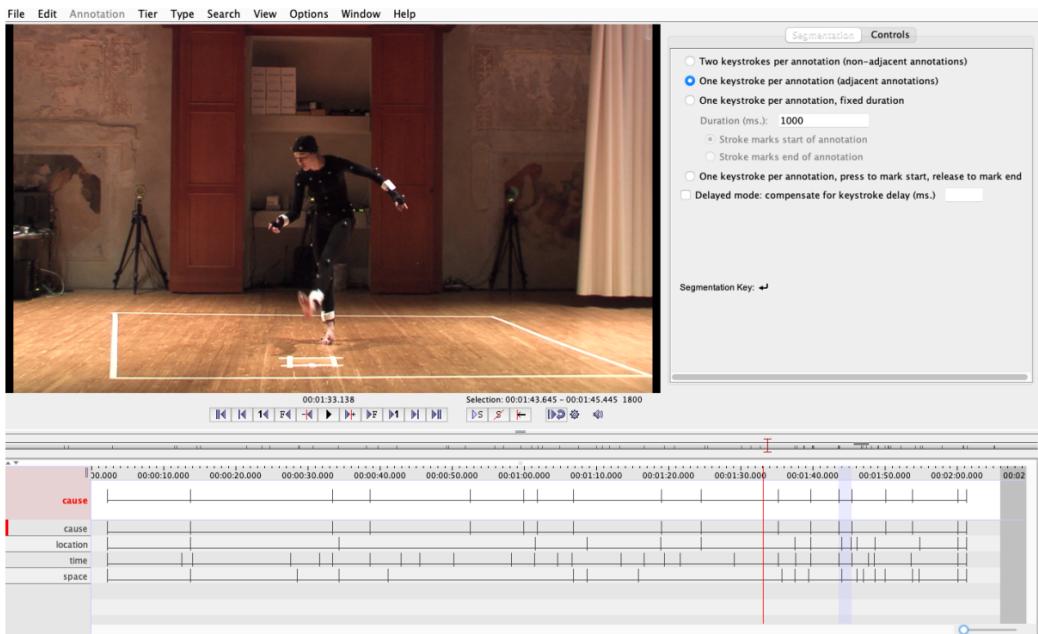


Figure 6: Annotation of the different salience with ELAN.

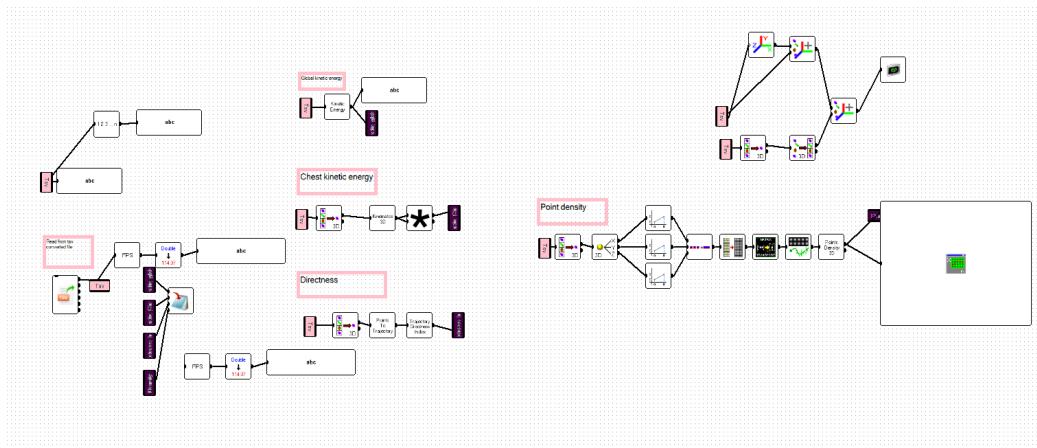


Figure 7: One of the patches implemented on EyesWeb for extracting the features.

2.3 Structure of workflow

In this chapter it will be explained which is the workflow that we did during our study (figure 8).

From the raw data analyzing to the salience classification and experiment testing:

1. Physical Layer: this layer was already achieved since the Motion Capture was given by the researcher at Casa Paganini - InfoMus [7] (chapter 2.2).
2. Salience and Ground Truth: After reviewing the material provided at the previous step, we evaluate which is the better definition of salience in our context, because we had to take in account the different possible movements of the dancers.
Then, we wrote down the Ground Truth based on that definition of salience (chapter 3).
3. *Takes* selection: With respect to the point 2 we had to choose wisely the *takes* which better represent the context, because not all the videos were adapt for representing our domain (chapter 5.1).
4. Data Cleaning: We cleaned up the data with a Linear interpolation (chapter 5.1.1).
5. Low Level Features: Then, we have defined the *features* at a lower level, global kinetic energy, Repetitiveness etc.
Hence these features have been collected thanks to EyesWeb (chapter 4).
6. Model Selection: We had explored two different approach for addressing the problem, a machine learning approach, and a Change Point Detection (CPD) technique (chapter 5 or 6).
7. Sliding Windows Template: For the specific model we built a template of sliding window able to convey the right piece of information basing on the used model (chapter 5.1.2 and 6.1).

8. Mid Level Features: For the Machine Learning (ML) approach, mid-level features have been selected, such as mean, variance, entropy etc. (Chapter 5.1.3)
9. Samples Definition and Downsampling: then for the ML model a specific definition of samples has been set (chapter 5.2).
10. Data Normalization: Data has been normalized (chapter 5.1.4)
11. Algorithm Selection: For the specific approach, a specific algorithm was used, for Machine Learning (ML), a Random Forest was employed for doing the classification (chapter 5.3).
12. Parameters Definition: Basing on the approach used parameters have been chosen, for ML a cross-validation has been done (chapter 5.3).
13. Model Evaluation: For finishing, different level of generalization has been explored (chapter 5.4).

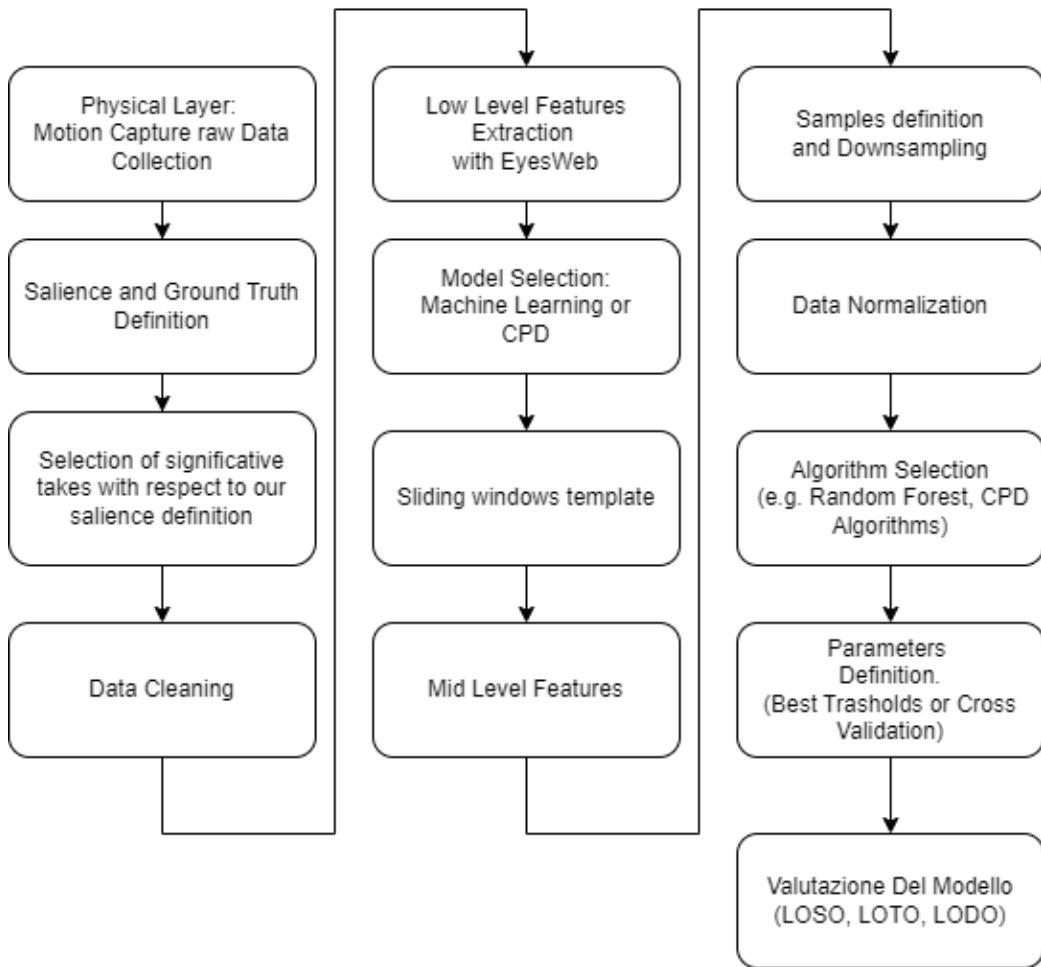


Figure 8: Workflow chart

3 Salience definition

When we started to approach this problem, one of the first challenge that we had to deal with, was to define what is a specific event that change in some way the behaviour of the dancer.

We address this event as a *Salience*.

The salience in the context of movement can be viewed in different ways, from an holistic one, to a granular one, by considering a movement within a large temporal window, or by viewing it with a short temporal span.

The problem can be addressed from different point of view, for instance, you can make a comparison between a movement which is more *impulsive* or more *sudden* [4], or you can look at the movement in terms of *Lightness* or *Fragility* [5].

The first issue that we have noticed is that, if we thought the salience as something which is too much related to the causes, and the psychological aspect of the human behaviour, any algorithm over that kind of approach gave us bad results, because addressing an intention or an emotion is a tricky challenge.

So, we have focused on looking more the movement of the dancer from the point of view of the Kinesiology, so if there is a movement behaviour which is similar with the observation of the whole set of frames which precedes that specific frame, then for us is not a salience, if instead the frame subjected to evaluation is the first frame before a whole set of frame which has different characteristics in terms of movement with respect to the set of frames before the specific evaluated frame, then for us it is a salience.

After that, by looking the set of salience resulting on the step before, and since the quantity of frames labeled "salience" is significantly lower with respect to the frame which are labeled "non-salience", we became aware that for an algorithm of machine learning this scenario could be a problem, because unbalanced dataset are often hard to address [10]. For trying to reduce the oddness between the two classes we had defined the concept of event as something which is more frequent, by using a more granular segmentation, hence visible changes in the movement such as an arm which starts to move with respect to the other one, or a repetitive movement and then a chaotic one, and so on...

3.1 Ground Truth

Then we thought to rebuild the Ground Truth, so with ELAN [1] we took annotations of all the events building upon the considerations made before (Figure 9). Since each event is detected between two sets of consecutive frames, events can be addressed as follows:

- S_1 - This event occurs when the dancer moves from a specific position in the stage to another position in the stage.
- S_2, S_3 - This event detect a change in the movement from the point of view of the repetitiveness, hence, if the dancer moves from a repetitive and orderly action, to a chaotic one (S_2), or vice versa (S_3).
- S_4, S_5 - It occurs when the dancer moves from a crouched position to an extended one(S_4), or vice versa(S_5).
- S_6, S_7 - Event occurred when the dancer moves from a static position of the head to a dynamic one (S_6), or vice versa (S_7).
- S_8 - The event which occurs when the dancer change the dynamic of her distal points: for instance if she moves her left wrist in the first set of frames, and then moves her right ankle, or if she moves her left ankle, and then moves her right wrist, or if she moves both the wrists and then stops to move one of the other distal parts.
Since there are 4 different distal parts, left wrist, right wrist, left ankle, right ankle, and all of them can be moving or not, we can think them as a set of boolean variables (1 if moving, 0 if not), we can think the vector $v_0=[0,0,0,0]$ as the vector which represents the situation in which the dancer does not move anything, and the vector $v_n=[1,1,1,1]$ in which the dancer moves all the distal parts (n=15).
Thus we have 2^4 different states (st_s) and hence the events (S_8) represent the transition from st_i to st_j ($i, j \in [0, 16], i \neq j$).

So, in total we have 23 different kind of salience.

The salience from S_1 to S_5 consider a more holistic change between the two set of frames, instead the set of salience which goes from S_6 to S_8 is referred

to a local movement.

During the phase of taking the Ground Truth, we did not make meaningful segments, we just consider the events occurrence, then for avoiding perceptual fusion problem (two different frame are perceived as one if the time between them is lower than 10 ms), we took the start of the segment (which is our salience) by looking frame per frame in the area of the perceived salience. For finishing, all the starting frames of the events have been exported in a text file where each value in that file represent the timestamp in seconds (with centesimal precision) of an occurrence of a frame representing a salience. For having the number of the frame instead of the timestamp we have computed the following formula: $n_f = \text{floor}(n_s * 50)$, since we have 50 frames per second (floor is the rounding down).

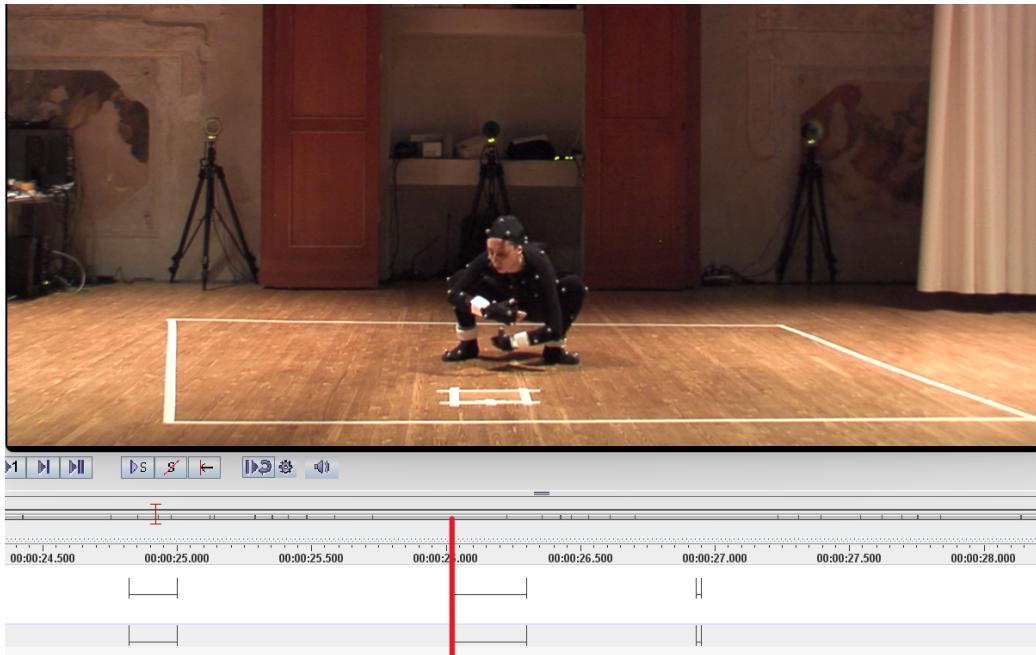


Figure 9:
The red line is over
the frame were
the salience is annotated (S_4).

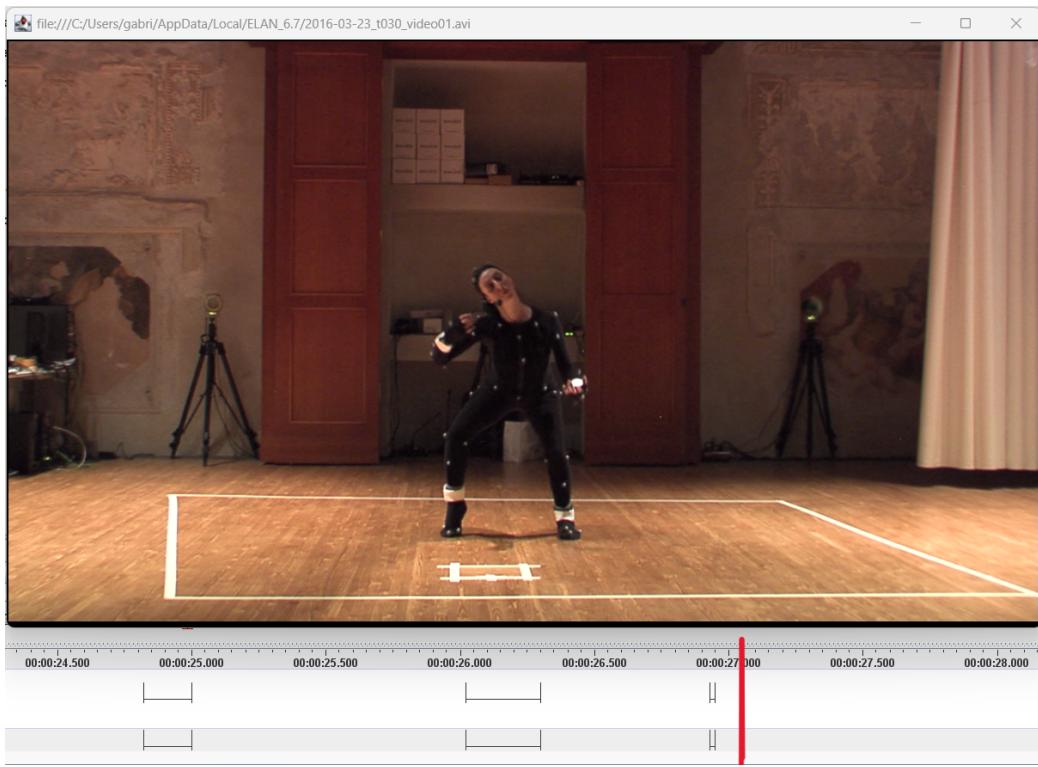


Figure 10:
After \approx 50 frames (1 second)
with respect to the
salience in Figure 9
the dancer is extended.

4 Low level features

In order to describe correctly the domain and to detect properly the different kind of salience, we had to find out which were the most representative features of our system.

Raw data was conceptually a three dimensional matrix in which on the first dimension we had the frame number, in the second dimension we had the ID_s of markers (e.g. ARIEL which refers to head marker, RPLM for the right wrist, and so on...), and on the third dimension we had the spatial axes (x,y,z), because each marker in each frame should have 3 values referred to its 3D location.

From these data we had to build significant low level features, because the spatial information of markers was not meaningful.

Since the salience were those listed in chapter 3, we focused on them for choosing the right features.

For this discussion we start to talk about the features from the locals (S_6, S_7, S_8) for arriving to the holistic ones ($S_1 - S_5$), because for addressing the more local salience, only one feature is enough, instead multiple features have been involved for a salience which is evaluated on the whole body like the S_1 .

Thus we prefer to define the feature of local salience as first.

All the features have been extracted by using an application (patch) developed for EyesWeb XMI, for each feature there will be an illustration on how the Objects within the application have been used.

Before talking about the effective low level features used in the process we want to mention the *Postural Tension* feature.

The Postural tension was a measure on how much the dancer's body is twisted, for evaluating that, unit vectors of head, shoulder, trunk and hips were compared in terms on how much they were parallel, the less they were parallel, the more the tension was higher.

At the beginning we thought that combining the postural tension to the others we could have good results, but by doing some tests we concluded that it was useless for our approach.

4.1 Distal Parts and Head movement - $S_8; S_6, S_7$

For representing the movement of the distal parts and the head, we considered the kinetic energy of the marker over the head and those related to the four different distal parts: right wrist, left wrist, right ankle, left, ankle. For doing that, we took from the TSV file the specific marker label:

- Left wrist = LPLM
- Right wrist = RPLM
- Left ankle = LHEL
- Right ankle = RHEL
- Head = ARIEL

In Figure 11 you can see the EyesWeb patch used for extracting the kinetic energy, this patch is used for all the markers listed before, by changing in the block employed for taking the 3D values of the markers, marker's label as input.

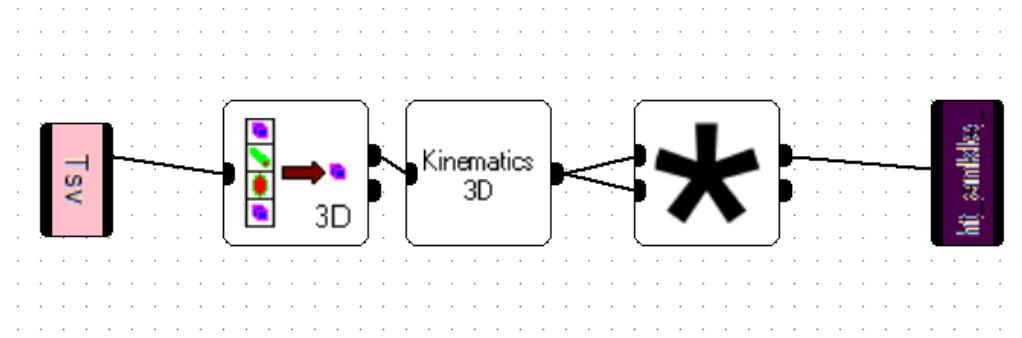


Figure 11: EyesWeb patch for extracting the kinetic energy from a marker

From the TSV all the markers values have been taken (pink block on the left), then the specific marker was chosen by using the 3D block, the Kinematics block extract the velocity of the markers along the time series, and then

thanks to the operator block (the fourth of the series), the square values have been calculated for having a good approximation of the kinetic energy by following the formula $KE = v^2$ (v stands for velocity). The last block store the results in a variable (purple block).

4.2 Point Density for crouched position - S_4, S_5

For measuring how much the dancer was crouched or extended we evaluated the *Point Density*, the concept behind this definition is the following: the more the distance among the different markers is higher on average (the dancer is extended), the more the value is higher, and vice versa (for the crouched position).

Even this value is extracted by using EyesWeb XMI (Figure 12).

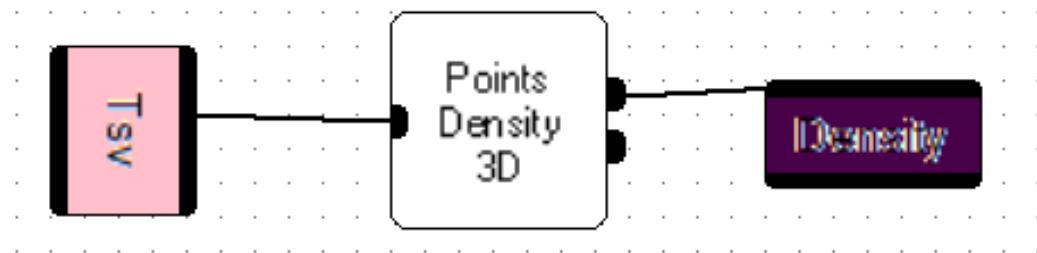


Figure 12: EyesWeb patch for extracting the Point Density

4.3 Repetitiveness - S_2, S_3

In order to have a measurement on how much the dancer make repetitive actions we took in consideration the *Repetitiveness* index of the Point Density, and kinetic energy of the distal parts.

Since we considered this feature as a mid level feature (we built it over low level features), it will be explained in the details on chapter 5.1.3.

4.4 Stage dancer position - S_1

Finishing, for the S_1 salience we thought as first to use a measure of general kinetic energy (Figure 13), because if a dancer moves toward the stage, of course the global kinetic energy, which is the sum of all the kinetic energies of all the markers, will increase, further a measurement of this value is enhanced by the kinetic energy of the ankles, because when the dancer is moving the kinetic energy of the ankles increase.

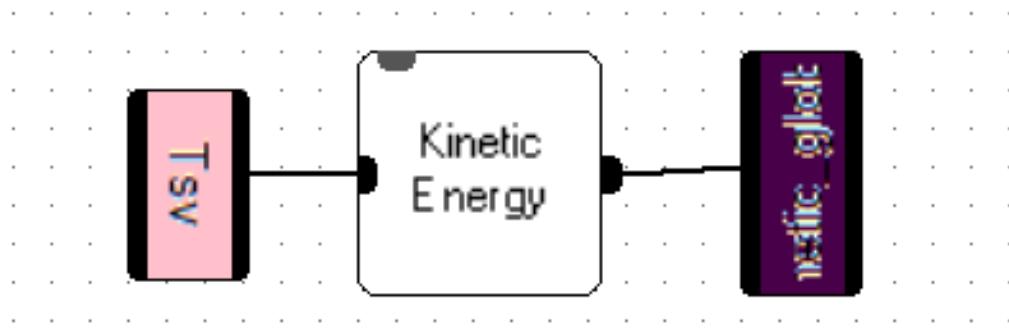


Figure 13: EyesWeb patch for extracting the Global kinetic energy.

4.5 Feature Vector

All the features stored in the EyesWeb XMI variables have been collected in a text file (Figure 14), hence, for each frame of the *take*, a row vector of low level features is build as follows: $F=[f_1, f_2, f_3, f_4, f_5, f_6, f_7, f_8]$.

According to the discussion made before, the several features are the following:

- f_1 = Global kinetic energy.
- f_2 = Postural Tension.
- f_3 = Head kinetic energy.

- f4 = Point Density.
- f5 = left wrist kinetic energy.
- f6 = right wrist kinetic energy.
- f7 = left ankle kinetic energy.
- f8 = right ankle kinetic energy.

The text file has a quantity of rows equal to the quantity of frame of the specific *take* and 8 columns, one for each feature.

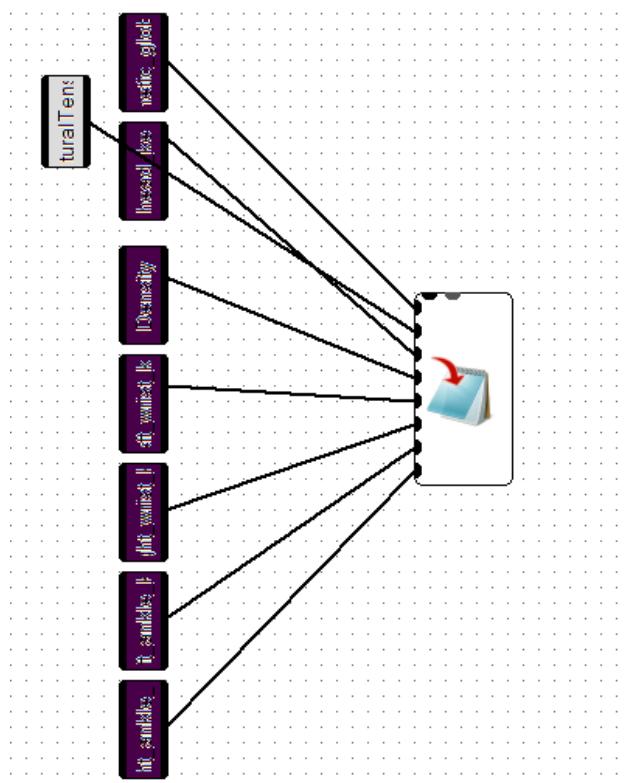


Figure 14: Saving the variables which store the features values over the frames on a text file (EyesWeb).

5 Supervised ML approach

Among all the existing methodologies for detecting a salience, we thought to address the problem by using a Machine Learning (ML) tool, of course one may think to use a CPD methodology or more in general, a statistic oriented algorithm, but since we had the possibilities to make a Ground Truth, and hence to use a ML algorithm which exploits its information (*Supervised*), we go throw that path.

Machine Learning is a branch of artificial intelligence focused on developing algorithms and models that allow computers to learn from data and make predictions or decisions without being explicitly programmed. It involves the utilization of statistical techniques to enable machines to improve their performance on a task through experience. In essence, Machine Learning enables computers to recognize patterns, infer insights, and adapt their behavior based on the data they are exposed to, thereby enabling them to solve complex problems and make smart decisions autonomously.

The next discussion is a brief introduction to ML technologies, hence analytical considerations behind the concept exposed will be omitted.

In general, the goal of a machine algorithm is to minimize a function which express the distance from the real output $y \in \mathbf{y}$, \mathbf{y} vector of outputs, with the estimated one \hat{y} .

For doing that, it learns patterns within a set of samples, which are called *Training Set*. During the *training phase* a machine learning algorithm learn which are the best *hyperparameters* (θ) and the best vector of weights ω that minimize the generic function:

$$L(y, \hat{y}).$$

And the model as well,

$$M(L(y, \hat{y}), \theta).$$

with L a function which express the error (distance) between y and \hat{y} , $\hat{y} = f(\mathbf{X}, \omega)$, \mathbf{X} set of samples and ω vector of weights.

Each sample can be viewed as a set of values belonging attributes called *features*.

For making an example, by imagining to use a ML model for predicting which is the age (\mathbf{y}) of a specific person, there can be defined a sample as a set of information about a person, for instance: height, weight, gender, etc.

Thus all these information are the features which can represent a person, and "80kg", "1.7 meters", "Female", etc. Is an instance of the set of features, and hence a *sample*.

So, having \mathbf{d} number of features, and \mathbf{n} number of samples, $X^{n,d}$ is our *Dataset*.

A feature can be in general categorical, or ordinal, in the first case values that have not a concept of distance among them, and hence that cannot be ordered (e.g. the gender in the example, or by remaining on a human dataset, the color of eyes), the second case, instead, has values which can be ordered (e.g. in the example, the weight or the height).

For how concern the output \mathbf{y} , it can be ordinal ($\mathbf{y} \in \mathbb{R}$), or it can be categorical as well (hence if \mathbf{C} is a set of categories, then $\mathbf{y} \in \mathbf{C}$).

If \mathbf{y} is categorical, then we talk about a *Classification* problem, if $|\mathbf{C}|=2$, we can call it *Binary Classification Problem*.

In this case there exists a lot of loss function (L) which can be used, among them, one of the most employed is the Categorical Cross-Entropy Loss:

$$CCEL = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{i,j} \log(\hat{y}_{i,j}).$$

Where N is the number of samples.

If \mathbf{y} is ordinal, we address the problem as a *Regression Problem*, since here there is a concept of distance between the predicted output and the actual one.

There can be employed for instance one of the following Loss functions:

- Mean Squared Error (MSE):

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2.$$

- Mean Absolute Error (MAE):

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|.$$

- Mean Squared Logarithmic Error (MAE):

$$MSLE = \frac{1}{N} \sum_{i=1}^N (\log(y_i + 1) - \log(\hat{y}_i + 1))^2.$$

If \mathbf{y} is known, the algorithm used will be in the context of the so called *Supervised* algorithms, if \mathbf{y} instead is unknown, the algorithm will be *Unsupervised*, in that last situation, the model will try to divide the samples in group (classes) and creating boundaries by observing the data itself.

The function $f(\mathbf{X}, \boldsymbol{\omega})$ is built during the training phase, given the set of samples \mathbf{X} (from which it tries to learn the patterns), it sets the weights $\boldsymbol{\omega}$ associated to the several features (it is a measure on which are the most important features for the sake of prediction for the specific dataset).

The *hyperparameters* θ are useful for avoiding the so called *Over-fitting*, this issue comes out when the function f fits too well the samples in the training set and then it doesn't generalize properly the dataset, intuitively this phenomena is recognized when the model predicts well the samples in the neighbor (with respect to the space of features) to those used in the training set, but it misses the classification with samples which are different.

For trying to find a properly good level of generalization (and so, if it's necessary reduce the complexity of function f), a phase called cross-validation is computed, the dataset is divided in two different sets, one for learning the model and one for testing the prediction quality for each "instantiation" of *hyperparameters* θ , once θ are found, the best vector $\boldsymbol{\omega}$ is defined as well and the function is built as follows:

Let $\boldsymbol{\omega}^*$ the best weights for the specific dataset, the best function is

$$f^* = f(\mathbf{X}, \boldsymbol{\omega}^*).$$

Below an example of function f :

$$f = \boldsymbol{\omega}^T \mathbf{X} + b$$

b is a bias (it can be 0). Hence a Machine Learning model can be:

$$M = \|\boldsymbol{\omega}^T \mathbf{X} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\omega}\|_2^2$$

This is the Regularized least square (RLS), because uses a least square as loss function, and it uses a regularizer (for handling the level of complexity) which is the term $\lambda \|\boldsymbol{\omega}\|_2^2$, λ is the hyperparameter (θ) and the subscript "2" in both terms stands for the norm L_2 . This is a model for regression problem, since f is linear it is used for problem which are not too complex.

The best model is given by searching $\boldsymbol{\omega}$ vector which minimize M by considering the right level of generalization (λ).

5.1 Data

In Our Context Data is given by Motion Capture (MOCAP) technologies, and stored in Tab-Separated Values (TSV) files (Chapter 2.2).

However, not all the *takes* were useful for our goal, so, before proceeding in cleaning the data we had to select them.

There were three different *dancers*: Cora Gasparotti, Muriel Romero and Marianne Gubri.

Each dancer has a specific number of *takes*, a *take* is a short video (50 fps, max 3 minutes) in which the dancer performed precise tasks for highlighting some specific movement, then by thinking at our salience definition, we have to do a selection.

For Cora Gasparotti 5 *takes* were chosen:

- t001: this *take* has captured our interest because she performed at the beginning a strike of fast and angular gesture, and on the second part of the video she was more smooth in the movement, then we had a lot of salience in the first part, and good frames which represent properly the non salience, in the second half.
This *take* has been annotated till 85 seconds (4250 frames).
- t004 (2019/08/08) - Even for this *take*, the gestures were very frenetic and angular on a first part, then the gesture has became less frenetic, till to became completely smooth in the end of the video.
Annotated 43 seconds (2150 frames).
- t004 (2019/05/22) - This *take* was easy to segment, hence all the salience were easy to see.
Annotated for 85 seconds (4250 frames).
- t005 - We thought that this was one of the most valuable *take*, all the movement were easy to segment, majorly for annotating the repetitiveness (S_2, S_3). Indeed this video was the one which perform worst in Leave One take Out (LOTO) with a f1_score of 0.66 (Chapter 5.4.2) with respect other videos, because it was used into test set and not for training.
Annotated for 85 seconds (4250 frames).

- t015 - Also for this video, it was easy to segment it was a good material for all the kind salience defined, then we used it.
Annotated for 45 seconds (2250 frames).

For how concern Marianne Gubri, the following 12 takes were chosen:

- t007 - This video was good because salience like when the dancer was crouched or extended (S_4, S_5) were identified clearly, good also for S_8 (distal points movement).
Annotated for 62 seconds (3100 frames).
- t008 - This *take* was easy to segment because movement were divided in time properly, useful for all the salience.
Annotated for 39 seconds (1950 frames).
- t010 - Very smooth movement and easy to segment, good extension of the body for pointing out S_4, S_5 salience.
This *take* has been annotated till 36 seconds (1800 frames).
- t018 - As for the previous one, even this *take* was useful for annotating clearly the salience, since the movement were good defined.
Annotated for 50 seconds (2500 frames).
- t019 - This *take* was full of salience, and easy to segment.
Annotated till 50 seconds (2500 frames).
- t024 - A video which alternated long static segment of equilibrium and (non salience), with brief segments in which there were few salience and very highlighted, in those frames the dancer had lost her equilibrium.
Annotated for 69 seconds (3450 frames).
- t026 - Express a huge quantity of salience well defined.
Annotated for 23 seconds (1150 frames).
- t041 - Also for this *take* there were a lot of salience, less then t024 or t026.
Annotated for 31 seconds (1500 frames).
- t042 - It is similar to t024, but in this case, when a salience occurred was more related to an interruption of extended slow movement, instead of a loss of equilibrium.
Annotated for 73 seconds (3650 frames).

- t043 - This *take* had all the salience in its frames, not a huge quantity and not defined perfectly as other videos, but good enough.
Annotated for 114 seconds (5700 frames).
- t047 - Full of salience well defined, the dancer did dynamic movement through all the *take* duration.
Annotated for 20 seconds (1000 frame).
- t048 - Similar to t047, but with less salience.
Annotated for 29 seconds (1450 frames).

For finishing, Muriel Romero's *takes* we had selected the following:

- t018 - This *take* had a large quantity of salience of almost all kind.
Annotated for 58 seconds (2900 frames).
- t026 - All kind of salience has been detected within this video, even the repetitiveness, it was a bit difficult to segment because of the large quantity of salience
Annotated for 145 seconds (7250 frames).
- t027 - Full of salience of all kind, difficult to making segmentation (as for t026 or t018). Annotated for 98 seconds (4900 frames).
- t030 - Very good for pointing out the salience for crouched and extended position (S_4 , $S - 5$).
Annotated till 81 seconds (4000 frames).

Hence we have approximately 66000 frames annotated, and according to the Ground Truth taken, only 641 were annotated as salience, thus we can say the dataset is extremely unbalanced.

As we previously said (In chapter 3), the choice of intending the salience even in a more granular point of view was useful for having an higher quantity of salience. In next chapters it will be explained also how this problem has been solved.

5.1.1 Data Cleaning

Before extracting from EyesWeb XMI the low level features by using the patch described in chapter 4, we had to deal with the problem of missing data (NaN), indeed, a lot of videos (*takes*) has been discarded because of the massive quantity of missing data, this phenomena is related to a wrong reception of the signal coming from the markers in the dancer's body (Chapter 2.2), often because markers are not visible from Qualisys Cameras [8].

After having discarded all the *takes* with too much Not a Number (NaN), we had to choose the policy for replacing them on the selected videos (which had a reasonable number of NaN).

There exists a lot of method for replacing the missing values, among them, the following [3]:

1. Mean/Median/Mode Imputation: Replace missing values with the mean, median, or mode of the non-missing values in the respective column. This is a simple and commonly used method but may not be suitable for all types of data.
2. Forward Fill or Backward Fill: For time-series data, fill missing values with the most recent non-missing value (forward fill) or the next non-missing value (backward fill).
3. Interpolation: Interpolate missing values based on the values of neighboring data points. Common interpolation methods include linear interpolation, polynomial interpolation, and spline interpolation.
4. Using Predictive Models: Train a machine learning model to predict missing values based on other features in the dataset. This approach can be effective but requires more computational resources and may introduce bias if the predictive model is not well-trained.
5. K-Nearest Neighbors (KNN) Imputation: Replace missing values with the average of the nearest neighbors' values. This method considers the similarity between data points and is suitable for datasets with continuous features.
6. Delete Rows or Columns: If missing values are too numerous or occur in specific rows or columns, you may choose to delete those rows

or columns altogether. However, this approach should be used with caution as it may lead to loss of valuable information.

7. Domain-specific Methods: In some cases, domain-specific knowledge can help in imputing missing values. For example, in financial data, missing values in certain columns could be replaced with zeros if it's reasonable to assume that missing values represent the absence of transactions.

For our scenario, since our data are continuous in time, and since *Nan*s often is a multiple occurrence (when markers are covered, multiple consecutive frames have been affected by *Nan*), we can use, for having a good approximation of the true value of the marker's position, a forward fill or a backward fill, or we can use an interpolation between the values before the missing values and after them.

Then we did a forward fill for the markers which starts from the frame 0 with *Nan* or a backward fill for missing data on the last frame of the specific *take*. For missing values which are in the middle of known values, we rather opted for a linear interpolation, by proceeding as follows: given a column (marker's axis position) we took the value known before the first *Nan*, let's say "fkb", and we took the first known after the last *Nan*, "fka", then we counted the quantity of unknown from the "fkb" to the "fka" indexes, we call it *n_nan*, then the computed expression for replacing the specific *Nan* value is:

$$\begin{aligned} val_diff &= fka - fkb \\ incr &= val_diff/n_nan \\ data[z, j] &= (round(fka + incr * (z - i + 1))) \end{aligned}$$

where *i* is the index of first *Nan*, *z* goes from *i* to the index of the last *Nan* and *j* is the column related to the specific axis (x,y,z) of the specific marker (round is a function for making the rounding, for instance to the third decimal place).

The cleaning from missing values can be improved for future utilization by using an interpolation which takes in account a group of frames after and before, instead of just looking just a frame (as linear interpolation does).

Another algorithm can consider the a spatial neighbor, instead of looking the frames before or after, it can look at the values taken by markers which are spatially closest to the missing one (KNN) and, after having evaluating

which are the closest markers and the values belonging to them, an estimation of the missing value can be computed.

5.1.2 Sliding Windows

5.1.3 Mid level features

5.1.4 Normalization

5.2 Samples Creation

5.3 Model Selection

5.4 Results

5.4.1 LOSO

5.4.2 LOTO

5.4.3 LODO

6 Statistic approach

6.1 Sliding windows

6.2 Mid Level Features

6.3 model

6.4 Results

7 Conclusions

7.1 Future Researches

References

- [1] The Language Archive. URL: <https://archive.mpi.nl/tla/elan>.
- [2] Ceccaldi E. “CEST: a Cognitive Event based Semi-automatic Technique for behavior segmentation”. PhD thesis. Università di Genova, 2022.
- [3] Roderick J Little and Donald B. Rubin. *Statistical analysis with missing data*. John Wiley & Sons, 2019.
- [4] R. Niewiadomski et al. “Automated Detection of Impulsive Movements in HCI”. In: *Conference CHItaly* (2015).
- [5] R. Niewiadomski et al. “Does embodied training improve the recognition of mid-level expressive movement qualities sonification?” In: *Multimodal User Interfaces* (2018).
- [6] *Oxford Latin Dictionary*. Oxford: Clarendon Press, 1982.
- [7] Casa Paganini. URL: <http://www.casapaganini.org/>.
- [8] Qualisys. URL: <https://www.qualisys.com/>.
- [9] *The Oxford English Dictionary*. Oxford: Clarendon Press, 1989.
- [10] L. Wang et al. “Review of Classification Methods on Unbalanced Data Sets”. In: *IEEE Access* (2021).
- [11] J.M. Zacks and K.M. Swallow. “Event segmentation”. In: *Current directions in psychological science* 16.2 (2007), pp. 80–84.