

Probability-based Dynamic Time Warping for Gesture Recognition on RGB-D data

Miguel Ángel Bautista^{1,2}, Antonio Hernández-Vela^{1,2}, Victor Ponce^{1,2}, Xavier Perez-Sala^{2,3},
Xavier Baró^{2,4}, Oriol Pujol^{1,2}, Cecilio Angulo³, and Sergio Escalera^{1,2}

¹*Dept. Matemàtica Aplicada i Anàlisi, Universitat de Barcelona, Gran Via 585, 08007
Barcelona, Spain.*

²*Centre de Visió per Computador, Campus UAB, Edifici O, 08193 Bellaterra, Barcelona, Spain.*

³*CETpD-UPC, Universitat Politècnica de Catalunya Neàpolis, Rambla de l'Exposició, 59-69,
08800 Vilanova i la Geltrú, Spain*

⁴*EIMT, Universitat Oberta de Catalunya, Rambla del Poblenou 156, 08018, Barcelona, Spain.*

Abstract. *Dynamic Time Warping (DTW)* is commonly used in gesture recognition tasks in order to tackle the temporal length variability of gestures. In the DTW framework, a set of gesture patterns are compared one by one to a maybe infinite test sequence, and a query gesture category is recognized if a warping cost inferior to a given value is found within the test sequence. Nevertheless, either taking one single sample per gesture category or a set of isolated samples may not encode the variability of such gesture category. In this paper, a probability-based DTW for gesture recognition is proposed. Different samples of the same gesture pattern obtained from RGB-Depth data are used to build a Gaussian-based probabilistic model of the gesture. Finally, the cost of DTW has been adapted accordingly to the new model. The proposed approach is tested in a challenging scenario, showing better performance of the probability-based DTW in comparison to state-of-the-art approaches for gesture recognition on RGB-D data.

Keywords: Depth maps, Gesture Recognition, Dynamic Time Warping, Statistical Pattern Recognition.

1 Introduction

Nowadays, human gesture recognition is one of the most challenging tasks in computer vision. Current methodologies have shown preliminary results on very simple scenarios, but they are still far from human performance. Due to the large number of potential applications involving human gesture recognition in fields like surveillance [3], sign language recognition [7], or in clinical assistance [4] among others, there is a large and active research community devoted to deal with this problem.

The recent irruption of easily available RGB-D sensors has reshaped the Computer Vision community allowing an extreme enrichment of visual data. This fact has enabled researchers to apply new techniques to obtain more discriminative features. As a consequence, new methodologies on gesture recognition can improve their performance by using RGB-D data.

From a learning point of view, the problem of human gesture recognition is an example of sequential learning. The main problem in this scenario comes from the fact that data sequences may have different temporal duration and even be composed of intrinsically a different set of component elements. There are two main approaches for this problem: On the one hand, methods such as Hidden Markov Models (HMM) or Conditional Random Fields (CRF) are commonly used to tackle the problem from a probabilistic point of view [7], especially for classification purposes. On the other hand, dynamic programming inspired algorithms can be used for both alignment and clustering of temporal series [8]. One of the most common dynamic programming methods used for gesture recognition is Dynamic Time Warping (DTW) [5].

However, the application of such methods to gesture recognition in complex scenarios becomes a hard task due to the high variability of environmental conditions. Common problems are: the wide range of human pose configurations, influence of background, continuity of human movements, spontaneity of humans actions, speed, appearance of unexpected objects, illumination changes, partial occlusions, or different points of view, just to mention a few. These effects can cause dramatic changes in the description of a certain gesture, generating a great intra-class variability. In this sense, since usual DTW is applied to compare a sequence and a single pattern, it fails when such variability is taken into account. We propose a probability-based extension of DTW method, able to perform an alignment between a sequence and a set of N pattern samples from the same gesture category. The variance caused by environmental factors is modeled using a Gaussian Mixture Model (GMM) [6]. Consequently, the distance metric used in the DTW framework is redefined in order to provide a probability-based measure. Results on a public and challenging computer vision dataset show a better performance of the proposed probability-based DTW in comparison to standard approaches.

The remaining of this paper is organized as follows: Section 2 presents the probability-based DTW method for gesture recognition, Section 3 presents the results and, finally, Section 4 concludes the paper.

2 Method

In this section we first describe the original DTW and its common extension to detect a certain pattern sequence given a continuous and maybe infinite data stream. Then, we extend the DTW in order to align several patterns, taking into account the variance of the training sequence by means of a Gaussian mixture model.

2.1 Dynamic Time Warping

The original DTW algorithm was defined to match temporal distortions between two models, finding an alignment/warping path between the two time series $Q = \{q_1, \dots, q_n\}$ and $C = \{c_1, \dots, c_m\}$. In order to align these two sequences, a $M_{m \times n}$ matrix is designed, where the position (i, j) of the matrix contains the alignment cost between c_i and q_j . Then, a warping path of length T is defined as a set of contiguous matrix elements, defining a mapping between C and Q : $W = \{w_1, \dots, w_T\}$, where w_i indexes

a position in the cost matrix. This warping path is typically subjected to several constraints:

Boundary conditions: $w_1 = (1, 1)$ and $w_T = (m, n)$.

Continuity and monotonicity: Given $w_{t-1} = (a', b')$, then $w_t = (a, b)$, $a - a' \leq 1$ and $b - b' \leq 1$, this condition forces the points in W to be monotonically spaced in time.

We are generally interested in the final warping path that, satisfying these conditions, minimizes the warping cost:

$$DTW(Q, C) = \min \left\{ \frac{1}{T} \sqrt{\sum_{t=1}^T M(w_t)} \right\}, \quad (1)$$

where T compensates the different lengths of the warping paths. This path can be found very efficiently using dynamic programming. The cost at a certain position $M(i, j)$ can be found as the composition of the Euclidean distance $d(i, j)$ between the feature vectors of the sequences c_i and q_j and the minimum cost of the adjacent elements of the cost matrix up to that point, i.e.:

$$M(i, j) = d(i, j) + \min\{M(i-1, j-1), M(i-1, j), M(i, j-1)\}.$$

Given the streaming nature of our problem, the input vector Q has no definite length and may contain several occurrences of the gesture pattern C . In order to detect the beginning and ending positions of the candidate gesture, the current ending cost is checked (the cost of the element in the last row). If this value is below a certain learned threshold μ , the warping path down to that matrix position is considered as a matching warping candidate gesture. An example of a begin-end gesture recognition together with the working path estimation is shown in Figure 3.

2.2 Handling variance with Probability-based DTW

Consider a training set of N sequences $\{S_1, S_2, \dots, S_N\}$ with each sequence S_g composed by a set of feature vectors at each time t , $S_g = \{s_1^g, \dots, s_{L_g}^g\}$ for a certain gesture category, where L_g is the length in frames of sequence S_g . Let us assume that sequences are ordered according to their length, so that $L_{g-1} \leq L_g \leq L_{g+1}, \forall g \in [2, \dots, N-1]$, the median length sequence is $\bar{S} = S_{\lceil \frac{N}{2} \rceil}$. This sequence is used as a reference, and the rest of sequences are aligned with it using the classical Dynamic Time Warping, in order to avoid the temporal deformations of different samples from a same gesture category. Therefore, after the alignment process, all sequences have length $L_{\lceil \frac{N}{2} \rceil}$. The set of warped sequences is $\{\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_N\}$. Once all samples are aligned, the feature vector corresponding to each sequence element at a certain time t among all warped sequences \tilde{s}_t is modelled by means of an G -component Gaussian Mixture Model (GMM) $\lambda_t = \{\alpha_k, \mu_k, \Sigma_k\}$, $k = 1, \dots, G$, α is the mixing value and μ and Σ are the parameters of each of the G Gaussian models in the mixture. As a result, each one of the GMMs that model each component of a gesture pattern \tilde{s}_t is defined as follows:

$$p(\tilde{s}_t) = \sum_{k=1}^G \alpha_k \cdot e^{-\frac{1}{2}(x-\mu_k)^T \cdot \Sigma_k^{-1} \cdot (x-\mu_k)}. \quad (2)$$

The resulting model is composed by the median sequence \bar{S} and a set of $L_{\lceil \frac{N}{2} \rceil}$ GMMs corresponding to the modeling of each one of the component elements of the warped sequence \tilde{s}_t for each gesture pattern.

2.3 Distance measures

In the classical DTW, a pattern and a sequence are aligned using a distance metric, such as the Euclidean distance. Since our gesture pattern is modelled by means of probabilistic models, if we want to use the principles of DTW, the distance needs to be redefined. In this paper we consider a soft-distance based on the probability of a point belonging to each one of the G components in the GMM, i.e., the posterior probability of x is obtained according to Eq.2. In addition, since $\sum_1^k \alpha_k = 1$, we can compute the probability of x belonging to the whole GMM as the following:

$$P(x, \lambda) = \sum_{k=1}^M \alpha_k \cdot P(x)_k, \quad (3)$$

which is the sum of the weighted probability of each component. An additional step is required since the standard DTW algorithm is conceived for distances instead of similarity measures. In this sense, we use a soft-distance based measure of the probability, which is defined as:

$$D(x, \lambda) = \exp^{-P(x, \lambda)}. \quad (4)$$

In conclusion, possible temporal deformations of the gesture category are taken into account by aligning the set of N gesture sample sequences. In addition, modelling with a GMM each of the elements which compose the resulting warped sequences, we obtain a methodology for gesture detection that is able to deal with multiple deformations in data. The algorithm that summarizes the use of the probability-based DTW to detect start-end of gesture categories is shown in Table 1. Figure 3 illustrates the application of the algorithm in a toy problem.

3 Experiments

In order to present the experiments, we discuss the data, methods and evaluation measurements.

3.1 Data

The data source used is the ChaLearn [2]¹ data set provided from the CVPR2012 Workshop challenge on Gesture Recognition. The data set consists of 50,000 gestures each one portraying a single user in front of a fixed camera. The images are captured by the KinectTM device providing both RGB and depth images. The data used (a subset of

¹ <http://gesture.chalearn.org/data/data-examples>

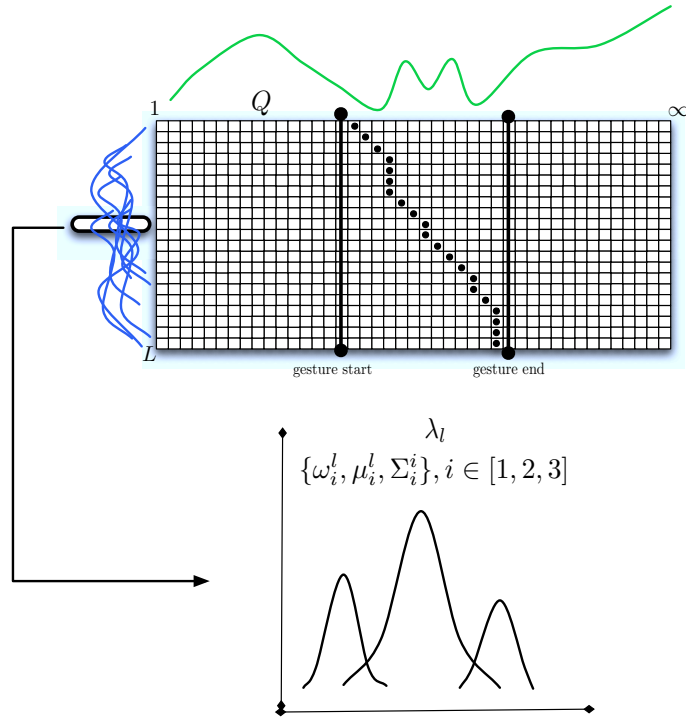


Fig. 1. Begin-end of gesture recognition of a gesture pattern in an infinite sequence Q using the probability-based DTW. Note that different samples of the same gesture category are modelled with a GMM and this model is used to provide a probability-based distance. In this sense, each cell of M will contain the accumulative D distance.

the whole) are 20 development batches with a manually tagged gesture segmentation. Each batch includes 100 recorded gestures, grouped in sequences of 1 to 5 gestures performed by the same user. For each sequence the actor performs a resting gesture between each gesture of the gestures to classify. For this data set, we performed background subtraction based on depth maps, and we defined a 10×10 grid approach to extract HOG+HOF feature descriptors per cell, which are finally concatenated in a full image (posture) descriptor. In this data set we will test the recognition of the resting gesture pattern, using 100 samples of the pattern in a ten-fold validation procedure. An example of the ChaLearn dataset is shown in Figure 2.

3.2 Methods and Evaluation

We compare the usual DTW and Hidden Markov Model (HMM) algorithms with our probability-based DTW approach using the proposed distance D shown in Equation 4. The evaluation measurements are the accuracy of the recognition and the overlapping for the resting gesture (in percentage). We consider that a gesture is correctly detected

Input: A gesture model $C = \{c_1, \dots, c_m\}$ with corresponding GMM models $\lambda = \{\lambda_1, \dots, \lambda_m\}$, its similarity threshold value μ , and the testing sequence $Q = \{q_1, \dots, q_\infty\}$. Cost matrix $M_{m \times \infty}$ is defined, where $N(x), x = (i, t)$ is the set of three upper-left neighbor locations of x in M .

Output: Working path W of the detected gesture, if any.

```

// Initialization
for  $i = 1 : m$  do
    for  $j = 1 : \infty$  do
end     $M(i, j) = \infty$  item    end
for  $j = 1 : \infty$  do
     $M(0, j) = 0$ 
end
for  $t = 0 : \infty$  do
    for  $i = 1 : m$  do
         $x = (i, t)$ 
         $M(x) = D(x, \lambda_i) + \min_{x' \in N(x)} M(x')$ 
    end
    if  $M(m, t) < \epsilon$  then
         $W = \{\text{argmin}_{x' \in N(x)} M(x')\}$ 
        return
    end
end
end

```

Table 1. Probability-based DTW applied to begin-end of gesture recognition.

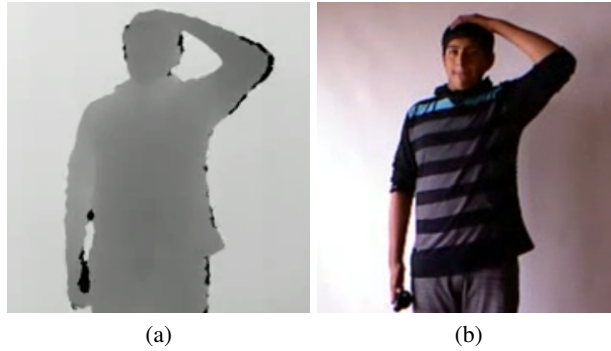


Fig. 2. Sample a) depth and b) RGB image for the ChaLearn database.

if the overlapping in the resting gesture sub-sequence is greater than 60% (the standard overlapping value [1]). The cost-threshold for all experiments was obtained by cross-validation on training data, and the confidence interval was computed with a two-tailed t-test. Each GMM in the probability-based DTW was fit with 4 components. For HMM, it was trained using Matlab toolbox Baum-Welch algorithm, and 3 states were experimentally set for the resting gesture, using a vocabulary of 60 symbols computed using K -means over the training data features. Final recognition is performed with

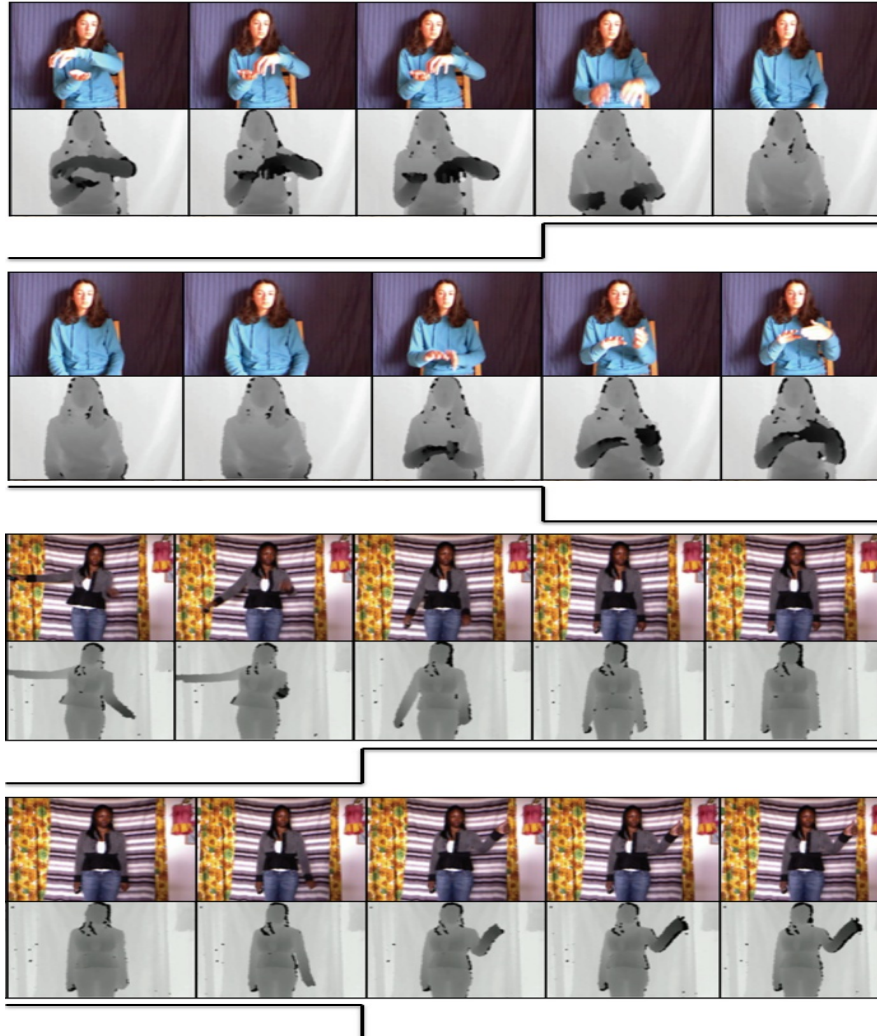


Fig. 3. Examples of resting gesture detection on the ChaLearn dataset using the probability-based DTW approach.

temporal sliding windows of different wide sizes, based on the training samples length variability.

Table 2 shows the results of HMM and the classical DTW algorithm, in comparison to our proposal on the ChaLearn dataset. We can see how the proposed probability-based DTW outperforms the usual DTW and HMM algorithms in both experiments. Moreover, confidence intervals of DTW and HMM do not intersect with the probability-based DTW in any case. From this results we can observe how performing dynamic

programming increases the generalization capability of the HMM approach, as well as a model defined by a set of GMMs outperforms the classical DTW on RGB-Depth data without increasing the computational complexity of the method.

Dataset	ChaLearn	
	Overlap.	Acc.
Probability-based DTW	39.08± 2.11	67.81±2.39
Euclidean DTW	30.03±3.02	60.43± 3.21
HMM	28.51±4.32	53.28±5.19

Table 2. Overlapping and Accuracy results of different gesture recognition approaches.

4 Conclusion

In this paper, we proposed a probability-based DTW for gesture recognition on RGB-D data, where the pattern model is learned from several samples of the same gesture category. Different sequences were used to build a Gaussian-based probabilistic model of the gesture whose possible deformations are implicitly encoded. In addition, a soft-distance based on the posterior probability of the GMM was defined. The novel approach has been successfully applied on a public RGB-D gestures dataset, being able to deal with multiple deformations in data, and showing performance improvements compared to the classical DTW and HMM approaches. In particular, the proposed method benefits from both the generalization capability from the probabilistic framework, when several observations of the training data are available, and the temporal warping capability from dynamic programming.

References

1. K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. Van Gool. A comparison of affine region detectors. *IJCV*, 65(1/2):43-72, 2005.
2. Chalearn gesture dataset, california, 2011.
3. A. Hampapur, L. Brown, J. Connell, A. Ekin, N. Haas, M. Lu, H. Merkl, and S. Pankanti. Smart video surveillance: exploring the concept of multiscale spatiotemporal tracking. *SPM, IEEE*, 22(2):38–51, 2005.
4. A. Pentland. Socially aware computation and communication. *Computer*, 38:33–40, 2005.
5. M. Reyes, G. Dominguez, and S. Escalera. Feature weighting in dynamic time warping for gesture recognition in depth data. *ICCV*, 2011.
6. M. Svensn and C. M. Bishop. Robust bayesian mixture modelling. *ESANN*, 64:235–252, 2005.
7. H.-D. Yang, S. Sclaroff, and S.-W. Lee. Sign language spotting with a threshold model based on conditional random fields. *IEEE TPAMI*, 31(7):1264–1277, 2009.
8. F. Zhou, F. D. la Torre, and J. K. Hodgins. Hierarchical aligned cluster analysis for temporal clustering of human motion. *IEEE TPAMI*, 2010.