

Trajectory-Based Recognition of Dynamic Persian Sign Language Using Hidden Markov Model

Saeideh Ghanbari Azar, Hadi Seyedarabi*

Faculty of Electrical and Computer Engineering, University of Tabriz, Tabriz, Iran

Abstract

Sign Language Recognition (SLR) is an important step in facilitating the communication among deaf people and the rest of society. Existing Persian sign language recognition systems are mainly restricted to static signs which are not very useful in everyday communications. In this study, a dynamic Persian sign language recognition system is presented. A collection of 1200 videos were captured from 12 individuals performing 20 dynamic signs with a simple white glove. The trajectory of the hands, along with hand shape information were extracted from each video using a simple region-growing technique. These time-varying trajectories were then modeled using Hidden Markov Model (HMM) with Gaussian probability density functions as observations. The performance of the system was evaluated in different experimental strategies. Signer-independent and signer-dependent experiments were performed on the proposed system and the average accuracy of 97.48% was obtained. The experimental results demonstrated that the performance of the system is independent of the subject and it can also perform excellently even with a limited number of training data.

Keywords: Sign Language Recognition, Persian Sign Language, Trajectory, Hidden Markov Model, Classification

1. Introduction

Sign language consists of a set of manual or non-manual gestures which are used for communication especially among deaf people. The majority of the gestures of a sign language

*Corresponding author

Email addresses: sghanbariazar@tabrizu.ac.ir (Saeideh Ghanbari Azar),
seyedarabi@tabrizu.ac.ir (Hadi Seyedarabi)

are manual gestures, while non-manual gestures like head movements (e.g. nodding), body movement (e.g. shrugging) and face expressions also play an important role in sign communications. Unfortunately, the use of sign language is usually restricted to the deaf community resulting in restricted communication of them with the rest of the world. Therefore, they are usually excluded from society and deprived of their rights to have equal educational and career opportunities. In order to address this problem, Sign Language Recognition (SLR) systems are developed to translate this language into speech or text. Developing efficient SLR systems can facilitate the communication of deaf people in society and remove the barriers for them.

One of the basic issues regarding SLR is that there is not a universal sign language. Sign languages of different countries have their own grammar rules. SLR systems have been rapidly developed in recent years for different sign languages including American [1–3], Chinese [4], Australian [5], Arabic [6, 7], Indian [8], Spanish [9] and Japanese [10]. For more reviews on sign language and different approaches developed for SLR systems refer to [11, 12].

Due to its broad range of capabilities, Machine Vision (MV) is the major tool used in the development of SLR systems. An MV-based SLR system usually consists of three components: hand tracker, feature extractor and classifier. Hand trackers job is to segment hand regions from the background of the input video frames. Some studies rely on data gloves to track the hand movements for hand tracking [6, 13–16]. Although these gloves are easy and precise to track, they usually contain heavy electromechanical devices which are inconvenient for the signers and limit their natural movements. Another type of studies relies on vision-based methods for hand tracking [7, 17, 18]. These techniques usually require some limitations on the signers cloths or imaging conditions. For instance, some of the vision-based studies need that the subject wear color gloves to facilitate the hand tracking process [19–21]. Nevertheless, these techniques are more convenient and cheaper. Feature extractor is the second stage of an SLR system. It takes the hand trackers data and produces a feature vector. Some SLR studies rely on hand shape information to extract a feature vector [21] while others rely on hand trajectory information [22]. Once feature vectors are

extracted, they need to be classified using an appropriate classifier. Many different classifiers have been utilized for recognizing different sign languages. These classifiers mainly include Neural Network (NN), K-Nearest Neighbor (KNN) and Hidden Markov Models (HMMs).

Developing machine vision based SLR systems started by the pioneering work of Starner et al. [23] in which they developed an American sign language recognition system. They placed the camera on top of a desk or a cap worn by the signer. They used two colored gloves to facilitate the hand tracking stage and classified the signs using HMM. In a recent study, Holden et al. [5] presented an HMM-based system which relies on hand shape information to extract a feature vector. This shape information includes hand size, direction, roundedness and the angle between 2 hands. Their system recognizes Australian sign language with the accuracy rate of 97% at the sentence level and 99% at the word level. Recently, many researches have been developed regarding Arabic sign language recognition. Al-Rousan et al.[18] suggested a vision-based system that uses Discrete Cosine Transform (DCT) and HMM for recognition of 30 Arabic signs in both offline and online modes. Recognition rates of 96.74% and 93.8% were obtained in offline and online mode, respectively. Tubaiz et al. [6] proposed a glove-based continuous SLR system. They used a feature extractor which emphasizes the temporal dependency of the data. A modified KNN approach is used for classification. The system recognizes 40 sentences of Arabic sign language with an accuracy of 98.9%.

Recently, deep learning-based approaches have proven to be very popular in many areas of machine vision applications including sign language recognition. The popularity of these approaches stems from their excellent discriminative abilities and successful performance. The seminal series of studies by Koller et al. [24–28] are among the first studies conducted for deep learning-based sign language recognition. In their early work [24], they have used deep learning for sign language recognition based on the shape of the mouth. In [27] they have developed an SLR system which uses Convolutional Neural Networks (CNNs) in an HMM framework. By combining the discriminative abilities of CNNs with the dynamic modeling ability of HMM, they have significantly improved the recognition performance on three benchmark sign language datasets, namely PHOENIX 2012 [29], PHOENIX 2014 [29],

and SIGNUM [30]. In their recent work [28], they have combined their previous works by adding multi-stream HMM to jointly solve the two sub-problems of hand gesture and mouth shape recognition. They develop a powerful and deep CNN with two bidirectional Long Short-Term Memory (LSTM) layers for recognition of continuous sign language sequences with weak and noisy labels.

1.1. Related Works to Persian Sign Language (PSL) Recognition

This part focuses on previous attempts made in recognition of PSL. Similar to other sign languages, PSL signs are divided into two main categories, i.e., static signs and dynamic signs. Static signs do not include hand movements and can be captured in a single image. Dynamic signs, on the other hand, include hand movements making it difficult to manipulate. Dynamic signs are usually captured in video frames, and video processing techniques are implemented to recognize them.

Development of a PSL recognition system is in its early stages. There have been few studies conducted in this field. These studies were all focused on image-based recognition of static signs. In the first PSL recognition system, Karami et al. [31] collected an image dataset of static alphabet signs. The images were then transformed into the wavelet domain. Different levels of wavelet transform including approximation coefficient of level 6, diagonal and horizontal details of level 6 and 7 and the vertical details of level 6 were used as feature vector. Multi-Layer Perceptron (MLP) neural network was used as a recognizer, and an accuracy rate of 94.06% was achieved. In another study, Barkoky and Charkari [32] designed a system for recognition of Persian sign numbers. They used a color-based technique to extract hand regions. A thinning method was then applied to these segmented images. The recognition was done by counting the number of endpoints of the thinned image. The accuracy rate of 96.6% was reported. In a similar study to [31], Moghaddam et al. [33] reported an image-based system to recognize alphabets of PSL. They used kernel based feature extraction methods including Kernel Principle Component Analysis (KPCA) and Kernel Discriminant Analysis (KDA). Three different classifiers including Minimum Distance (MD), Support Vector Machine (SVM) and NN were used to compare the results. In a more

recent study, Zare et al. [34] proposed a recognition system for 10 Persian static signs including six numbers and four words. They used skin segmentation in different color spaces to detect the hand regions. They employed Fourier descriptors as features to train a classifier. Their approach performs good results for real time signer-independent recognition system.

Each of the works discussed above for PSL recognition has its advantages. Since PSL recognition is a newly evolving field of study, there remains some challenges which motivates this paper. All the introduced PSL recognition systems are developed for alphabet or number signs which are not very useful in everyday conversations of the deaf community. This problem indicates the need for developing a dynamic sign recognition system which can recognize more practical signs. Motivated by this, we present a dynamic PSL recognition system. Over 1200 videos of 20 dynamic signs are collected for this system. A region growing technique is used to extract the motion trajectory of the hand and three other shape information. HMM with Gaussian observations is finally utilized to classify 20 dynamic signs.

To summarize, we make the following contributions: First, a new dynamic PSL dataset with 1200 videos is collected which contains 20 single-handed signs that are practical in everyday communication of the deaf community. To increase the diversity of the dataset, 12 individuals are participated in this dataset making it more subject-independent. Second, a dynamic PSL recognition system is proposed which uses a simple trajectory extraction approach based on region growing. The system performs excellently independent of the subject performing the signs, and it has the accuracy of more than 95% even with a limited number of training data.

The rest of the paper is organized as follows. Section 2 describes the collected dataset. Section 3 and Section 4 elaborate the trajectory extraction approach and the HMM classifier, respectively. Section 5, presents the experiments and compares the results. Finally, the paper is concluded in Section 6.

2. Dataset Collection

Since there is no dynamic PSL dataset available, the authors needed to construct a dataset. The dataset was collected in the Society of Deaf People (SDP), Urmia, Iran and

Table 1: List of the signs of the dataset.

Sign Code	Sign	Sign Code	Sign
1	Sad	11	Eat
2	Wish	12	Sun
3	Dear	13	Mother
4	Sorry	14	People
5	How?	15	Go
6	Student	16	Day
7	Today	17	Hear
8	Forget	18	Brave
9	Please	19	Natural
10	Danger	20	Can

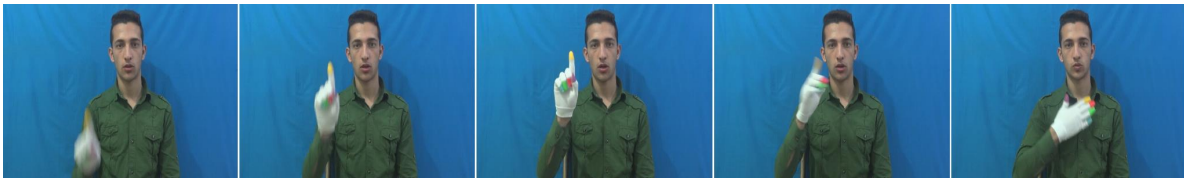


Figure 1: Samples of the captured video frames for sign wish.

named as University of Tabriz Persian Sign Language dataset (UoT-PSL) [35]. In order to develop an efficient recognition system, the collected dataset needs to contain different versions of performing a sign. This can decrease the dependency of the system on the subject. For this purpose, twelve volunteers including six male and six female signers were participated in this dataset. The signers included both deaf and hearing individuals. Twenty dynamic signs of PSL were chosen to be included in the dataset. The signs were selected from the single-handed signs appearing in everyday conversations with the consultation of the experts in SDP. A list of the signs available in this dataset is presented in Table 1. Each individual performed a sign 5 times producing 60 samples for each sign. In the following parts of this paper, the signs will be referred to by their corresponding code in Table 1.

A Sony digital camera (model DSC-HX9V) was used to obtain a total number of 1200 videos from 20 signs. There were no particular restrictions on the light of the environment. The collected videos were in AVI format and the frame rate was set to 25 frames per second and the spatial resolution was 1800×2000 pixels. The audio contents of the videos were eliminated to avoid unnecessary complexities.

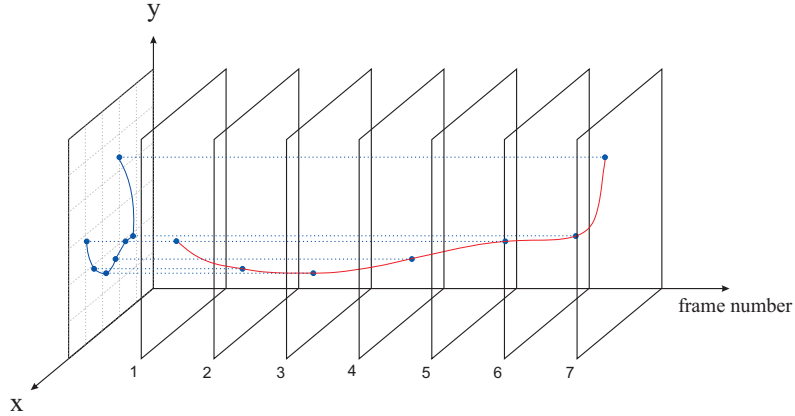


Figure 2: A representation of the hand trajectory.

To ease the process of hand tracking, some restrictions were imposed on the imaging conditions. The signers stood in front of a blue background and were asked to wear a simple white glove. Colored marks were considered on the fingertips of the glove for possible future studies but they were not utilized in the present study. All the signs were chosen from single-handed Persian signs to avoid possible occlusions of the hands. Figure 1 presents a sample of the captured video frames.

3. Sign Trajectory Extraction

The first step in developing the PSL recognition system is hand trajectory extraction. That is, in each frame, the hand region is detected, and its centroid in x and y-axes are saved. For a video stream these extracted centroids form the hand trajectory. Figure 2 gives an illustration of this definition. The hand region extraction procedure is explained in the following.

It should be noted that 10-15 starting frames of the videos did not contain the hand of the signer. In this paper the frame in which the hand appears for the first time in the scene is referred to as *start frame*. Thus, the first step was to detect this so-called start frame. There is no hand region in the frames that come before the start frame. Therefore, they all contain the same image of the signer. This means that subtracting these frames from the first frame results in an approximately zero image. This fact was exploited to detect the

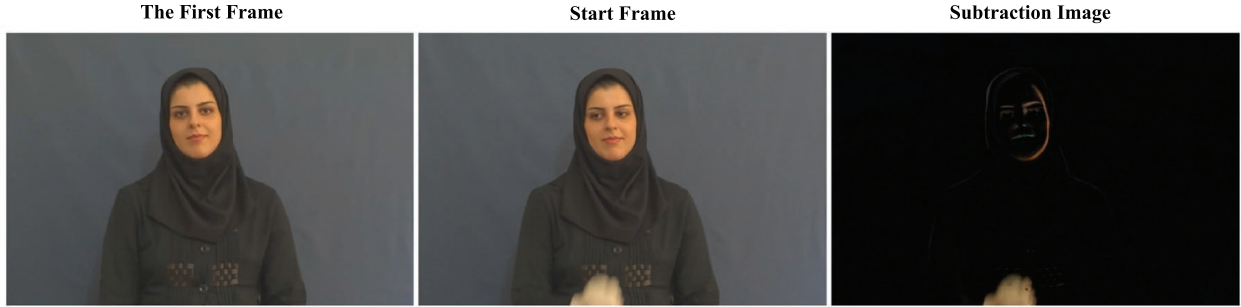


Figure 3: An example of the first frame of the video, start frame and their subtraction image

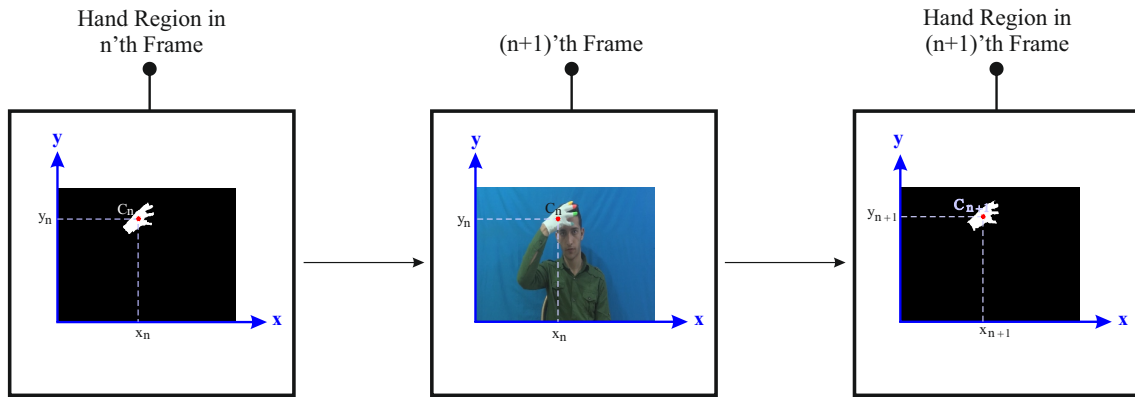


Figure 4: Hand region extraction of $(n + 1)$ th frame using the hand centroid of the n th frame as seed for region growing.

start frame. In this regard, each frame was subtracted from the first frame and among this subtracted stream of frames the first one containing a nonzero region bigger than a threshold determined the start frame. Figure 3 presents an example of the first frame of the video, start frame and their subtraction image.

After detecting the start frame, a hand tracking process begins to extract the trajectory of the hand. This hand tracking was accomplished via a region growing technique and is illustrated in Figure 4. Specifically, the centroid of the hand region in the start frame was obtained and was denoted as $C_1 = (x_1, y_1)$. This centroid's location was used as a seed for region growing in the frame that comes after the start frame. This region growing produces the hand region in this frame. Then, the hand region centroid of this frame denoted as

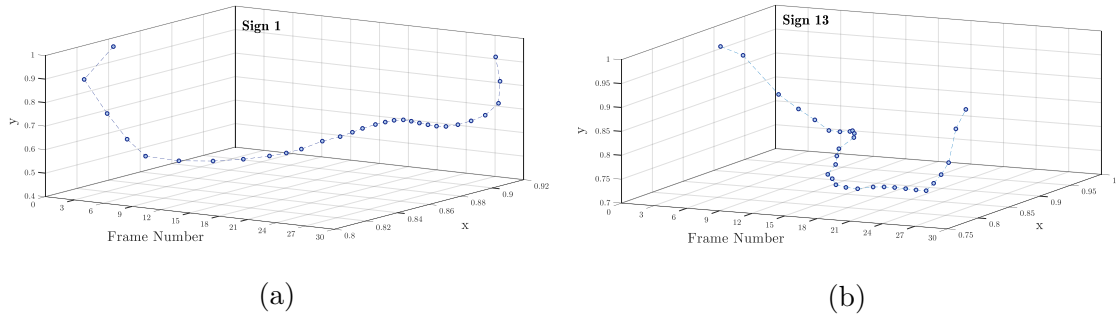


Figure 5: Two examples of the extracted trajectories for signs 1 (sad) and 13 (mother).

$C_2 = (x_2, y_2)$ was used as a seed for region growing in its next frame. Figure 4 illustrates the hand region extraction of the $(n + 1)$ th frame using the hand centroid of the n th frame, i.e. C_n . For all the frames of the video, this procedure was repeated, producing the hand trajectory. Figure 5 shows two examples of the extracted trajectories. In addition to the centroid of the hand, three simple shape information were also extracted from each frame as features. These shape features include: area, orientation and eccentricity. Area determines the number of pixels in the hand region for each frame. Considering the hand region as an ellipse, orientation measures the angle between its x-axis and major axis. Eccentricity is a measure of how much the bounded ellipse of the hand deviates from being circular. These features were added to the hand centroids forming a five-dimensional feature vector for each sign.

These time-varying feature vectors will be used as the observation sequences of the HMMs. Since the signs in the dataset were performed with different subjects and each subject had his/her own speed of performing the sign, there are vast differences in the number of the frames of the videos. To decrease the subject dependency of the system, we need to normalize the number of frames before training the HMMs. For this purpose, a linear temporal interpolation with 30 query points was used to normalize the number of frames. Therefore, for each video sample we extracted a 5×30 feature matrix.

4. Hidden Markov Model Based Classification

Unlike static signs which create time-invariant features, dynamic signs produce features which vary in time. In order to classify these time-varying features, we need a system that can model this dynamic nature of the features. HMM has long been used for the classification of temporal patterns [36], and it has been proved to be successful in sign language classification [23]. This section gives a brief introduction to HMM.

HMM is a stochastic model which contains a Markov chain with an invisible or hidden sequence of states. If we denote the number of hidden states as q , an HMM can succinctly be represented by:

$$\lambda = (\Pi, \mathbf{A}, \mathbf{B}) \tag{1}$$

where Π is a $q \times 1$ matrix containing the initial probabilities of the states and \mathbf{A} is the transition probability matrix. \mathbf{B} is called the state emission probability distribution and its components are denoted as b_j for j th state. At each time, the process is in one of the hidden states and generates observations according to these emission probability distributions. The observations either can be discrete or continuous. For discrete observations, the emissions of each state are represented by probability mass functions (pmf), and for continuous observations, they are represented by probability density functions (pdf). Refer to [37, 38] for detailed tutorial on HMM.

The observations used in this study are five-dimensional continuous feature vectors. Therefore, a pdf should be assigned for estimating these observations. The mixture of Gaussians is proved to be a successful method for estimating the pdf of continuous observations [39]. Hence, we model the observations of each state of the HMM with a mixture of Gaussians. Let the d -dimensional observation vector of each state be denoted as \mathbf{x} and the state at time t be denoted as q_t . The pdf of the observation at state j can be modeled as:

$$b_j(\mathbf{x}) = P(\mathbf{x}|q_t = j) = \sum_{i=1}^M c_i N(\mathbf{x}, \mu_i, \Sigma_i) \tag{2}$$

where M is the number of mixing Gaussian pdfs and c is the mixing parameter satisfying:

$$\sum_{i=1}^M c_i = 1. \tag{3}$$

$N(x, \mu, \Sigma)$ is a multivariate Gaussian distribution with corresponding mean vector μ and covariance matrix Σ .

HMM-based classification is performed in two steps, i.e. training and evaluation. In training step, an HMM is trained for each sign. That is, the parameters of the triplet λ are estimated. This procedure is known as the *training problem* of HMM. The parameters are initialized to random values and then estimated using the BaumWelch algorithm [37]. Assuming we have the total number of S classes or signs, the result of training step will be S trained HMMs represented as $\lambda = \{\lambda_1, \lambda_2, \dots, \lambda_S\}$. In the evaluation step, once HMMs are trained, a test sign with observation vector \mathbf{x} , is recognized by computing the probability of \mathbf{x} given each trained HMM. This is known as the *evaluation problem* of HMM and is solved using the forward-backward algorithm [37]. Specifically, given a set of trained HMMs for S signs, i.e. $\{\lambda_1, \lambda_2, \dots, \lambda_S\}$, and the observation sequence of the test sign, its sign label is assigned according to the following formula:

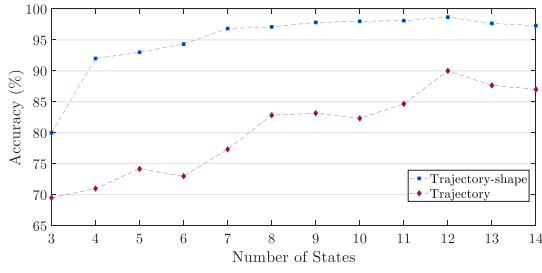
$$\operatorname{argmax}_{i=1,2,\dots,S} P(x|\lambda_i). \quad (4)$$

5. Results and Discussion

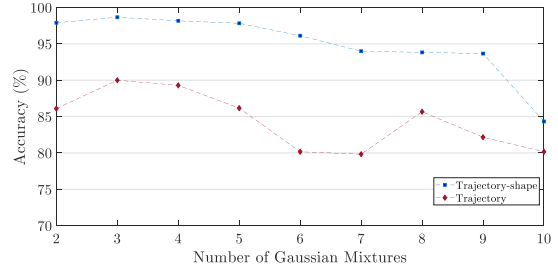
In this section, the performance of the proposed Persian sign language recognition system is evaluated in different experiments. Before conducting the experiments, the main parameters of the system are tuned. All the following experiments are performed with two sets of features. In the first set, which is referred to as *trajectory* features, only the x-y position of the hand is used as features. In the second set, which is referred to as *trajectory-shape* features, in addition to hand position, the shape information of the hand is also used as features. That is, the first set contains 2-dimensional features while the second one contains 5-dimensional features.

5.1. Parameter Tuning

To achieve best models for sign trajectories, there are two main parameters that need to be tuned, namely the number of hidden states and number of Gaussian mixtures. For

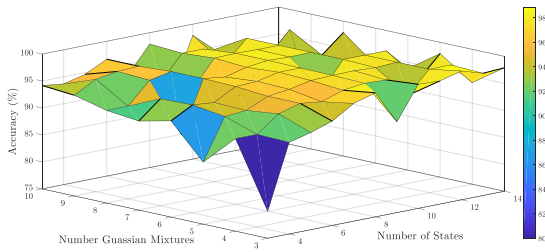


(a)

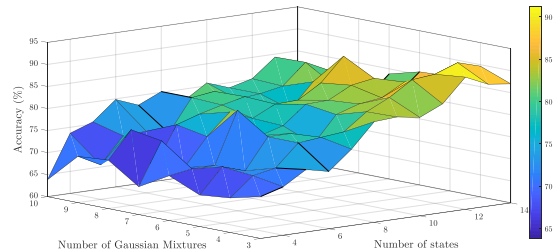


(b)

Figure 6: Tuning the number of states and the number of Gaussian mixtures. a) Classification accuracy as a function of number of states. b) Classification accuracy as a function of number of Gaussian mixtures.



(a)



(b)

Figure 7: Classification accuracy as a function of the number of states and number of Gaussian mixtures. a) For trajectory-shape features. b) For Trajectory features.

this purpose, 20% of the samples of each sign were randomly chosen for training and the rest of the data were used for testing. For an HMM, the most important parameter is the number of hidden states. Figure 6a shows the accuracy of the system for different number of states while fixing the other parameters. It can be observed that for both set of features, the accuracy of the system increases as we increase the number of states from 3 to 12. The highest accuracy of the system is achieved for 12 states for both sets of features, i.e., 98.66% for trajectory-shape and 91.2% for trajectory features.

The next parameter to be tuned is the number of Gaussian mixtures. Figure 6b shows the accuracy of the system for different number of Gaussian mixtures while fixing the other parameters. For both sets of features, the best accuracy is achieved for 3 mixtures, and it decreases as we increase the number of mixtures, revealing that 3 mixture of Gaussians is

the best representation for our observation data. To more evaluate the role of these two parameters, Figure 7 presents the classification performance for varying number of states and Gaussian mixtures. From this figures similar deduction to Figure 6 can be made about the optimal number of states and mixtures. Moreover, it can be observed from the figure that the trajectory-shape set of features (Figure 7a) is more robust to these parameters than the trajectory features (Figure 7b). To summarize, considering the results presented in Figure 6 and Figure 7, the optimal number of states and Gaussian mixtures were set to 12 and 3, respectively.

5.2. Sign Classification

In this section, the classification results of 20 dynamic Persian signs are presented. After extracting the hand trajectory and shape information, 20 HMMs were trained using both sets of features with 12 hidden states and 3 mixtures of Gaussian. In addition to HMM, Support Vector Machine (SVM) with polynomial kernel was also used for classification of signs and the results were compared to the ones obtained from HMM classification.

Three different training strategies with 20% of the samples of each sign were used to evaluate the performance of the system, namely random, subject-dependent and subject-independent training. In random training strategy, as its name suggests, 20% of the samples were selected randomly as training data, leaving the rest of the data for testing. In signer-dependent strategy, 20% of the samples of each subject were selected as training data. Considering we have five samples from each subject, one sample per subject was selected as training data. That is, we made sure that each subject had a sample among training data. In subject-independent strategy, on the other hand, we trained the system with samples of only two subjects, and the samples from the other ten subjects were left for testing. The results are presented in Table 2. All the experiments were conducted in 10 runs, and the mean and variance values of the classification accuracy are reported in Table 2.

Some observations can be made from this table. First, in both classifiers, adding shape information to the trajectory features has significantly (10%) increased the accuracy of the system, indicating the importance of the hand shape information in sign classification.

Table 2: Classification results for different classifiers with 20% of the samples used for training in different training strategy.

Classifier Feature set	SVM		HMM	
	Trajectory	Trajectory-shape	Trajectory	Trajectory-shape
Random	78.47 (± 0.85)	87.77 (± 1.12)	87.12 (± 0.21)	98.13 (± 0.11)
Subject-dependent	83.54 (± 0.62)	89.79 (± 0.23)	87.01 (± 0.13)	97.63 (± 0.12)
Subject-independent	62.90 (± 0.75)	67.10 (± 0.41)	83.20 (± 0.31)	96.70 (± 0.09)

Second, the results obtained by HMM is notably better than SVM. This may be due to the inability of SVM in dealing with the time-varying nature of the features, while HMM can successfully model these temporal features. Third, considering different training strategies, the results obtained by SVM meaningfully decrease in the subject-independent case and slightly increase for subject-dependent training. This leads to the conclusion that the system designed with SVM as a classifier is extremely subject-dependent. Contrary to SVM, HMM represents excellent subject-independent results and the classification accuracy drops only 1.4% in subject-independent case.

To further discuss the concept of signer-independency, Figure 8 presents samples of the same sign performed by a single signer (Figure 8a) and three different signers (Figure 8b). As it can be seen in Figure 8a, different realizations of a sign performed by a single signer are very similar in terms of both the shape of the trajectory and the x and y values of each frame number. For different signers (Figure 8b), although the x and y values of each frame are different, the shape of the trajectory is almost similar. SVM treats each frame as a separate feature and fails to see the temporal pattern of the features. As a result, it hardly recognizes the similarities between the signs performed by different signers. Therefore, exhibits weak signer-independent results. HMM, on the other hand, can model this temporal patterns using its state transition capabilities.

The trajectories of the 20 signs of the dataset are illustrated in Figure 9. This figure portrays the spread of realizations between these 20 signs. As it can be seen in this figure,

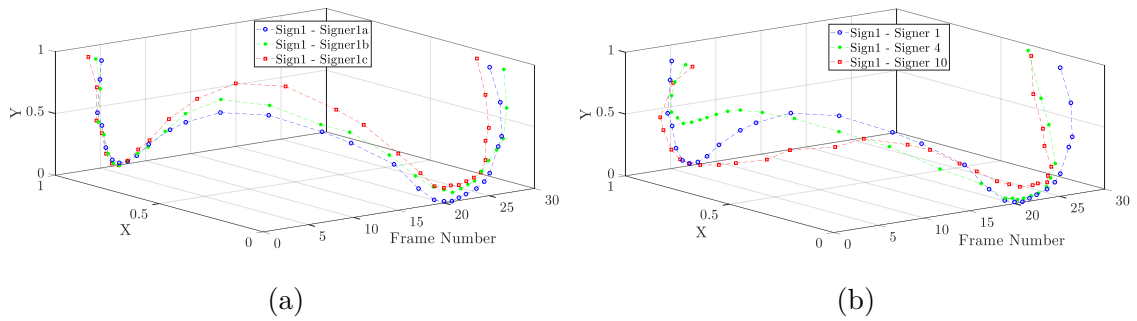


Figure 8: Samples of the same sign performed by (a) a single signer and (b) three different signers.

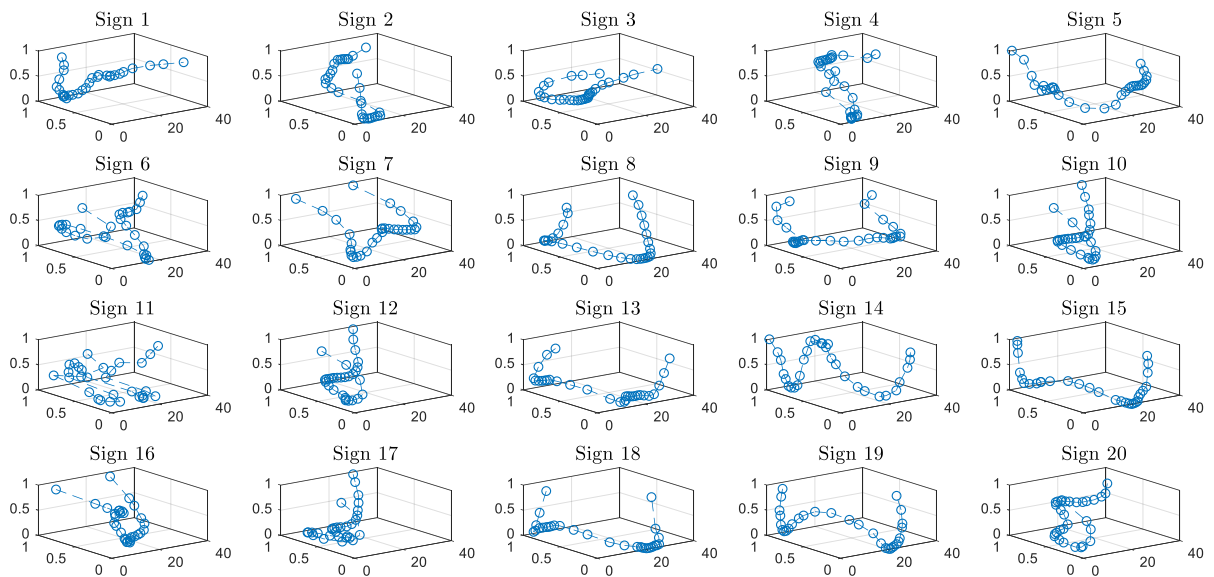


Figure 9: The trajectories of the 20 persian signs used in this study. Note that the labels are eliminated to decrease the ambiguity of the figure.

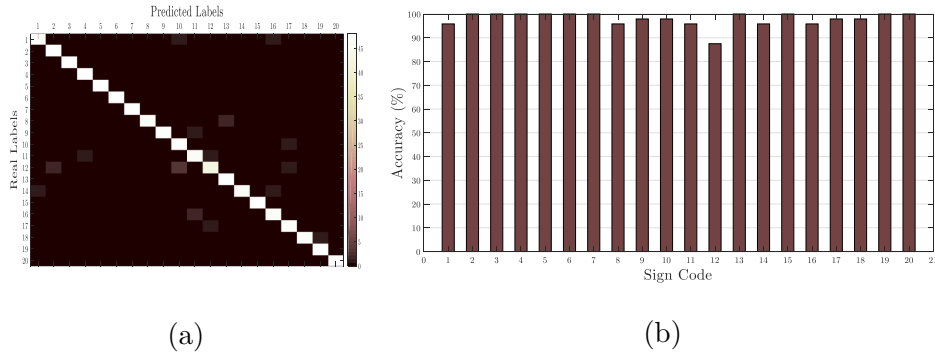


Figure 10: Performance of the proposed HMM-based system. Each sign is represented by its corresponding code. a) Confusion matrix of the classification. b) Accuracies obtained for each class of signs.

most of the signs are distinguishable while a few of them have similar trajectories that can challenge the performance of the system. To better evaluate the performance of the proposed HMM-based system, the confusion matrix of the classification and the accuracies obtained for each class of signs are illustrated in Figure 10a and Figure 10b, respectively. In these figures each sign is represented by its corresponding code from Table 1. According to these figures and the sign trajectories of Figure 9, following observations can be made. Half of the signs (signs 2-7, 13, 15, 19 and 20) are classified with 100% accuracy. Among these signs, signs 2-7 and sign 20 have distinguishable trajectories and the obtained accuracy was predictable. Signs 13, 15 and 19 have similar trajectories but they have been classified with 100% accuracy. This can be explained by the added shape information that has enable the system to discriminate between these signs. The lowest accuracy is obtained for sign 12 and it is mainly misclassified with sign 10, which may be due to their similar trajectories.

One of the most critical aspects of a recognition system is its level of dependence on the number of training data. Regarding the challenges in the acquisition of the sign videos, the number of available samples for each class is usually restricted. Therefore, it is essential for a recognition system to perform correctly with limited training data. Figure 11 exhibits the accuracy of the examined methods as a function of the train data percentage. It can be seen that the SVM-based methods rely significantly on the number of training data and their performance decrease as we decrease the train data, whereas HMM-based methods,

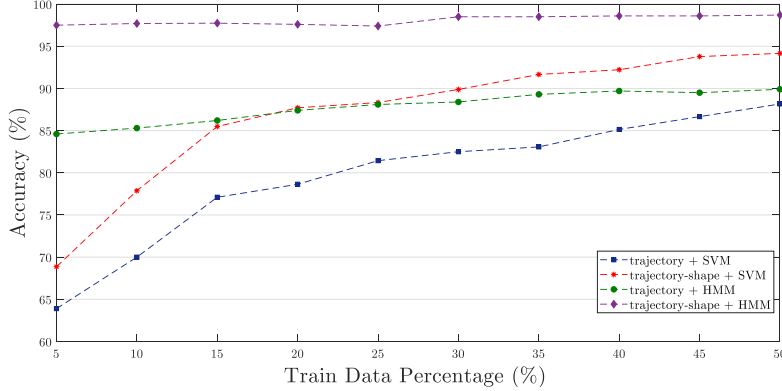


Figure 11: The accuracy of the examined methods as a function of the train data percentage.

especially the method with trajectory-shape features, are entirely robust to the number of train data. It can be observed from the figure that even with 5% of the data for training, the system can successfully model the signs, and for train data percentage of more than 30, the performance of the system remains almost the same. Note that the accuracy is not yet at ceiling for training data percentage of 50.

6. Conclusion

In this study, a dynamic Persian sign language recognition system is presented. A dataset containing 1200 videos of 20 signs were collected. Hand trajectories along with three hand shape information were extracted from video frames using a region growing technique. HMM with Gaussian mixture observations was utilized to model these trajectories and their temporal patterns. According to the experimental results, the HMM-based system with hands trajectory and shape information as features can successfully recognize these 20 signs with an average accuracy of 98.13%. Moreover, the experiments indicated that the performance of the system is independent of the subject, and it has excellent performance even with a limited number of training data.

This study being an initial study on dynamic PSL recognition has only focused on the trajectories of the signs. While, it is likely that using a wider dictionary of signs will increase the possibility of more similar trajectories leading to a need for more training data.

This problem can be addressed by using two cameras to extract both spatial and depth information and decrease the possibility of similar trajectories. Another solution may be to use more sophisticated approaches like deep learning based features. For future studies, the authors will be focused on updating the dataset and using deep learning based approaches for PSL recognition.

Acknowledgment

This paper is published as part of a research project supported by the University of Tabriz, Research Affairs Office, Iran. The authors would like to thank the Society of Deaf People (SDP), Urmia, Iran, for the many valuable assistances they provided during the acquisition of the dataset.

References

- [1] C. Vogler, D. Metaxas, Parallel hidden markov models for american sign language recognition, in: *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*, Vol. 1, IEEE, 1999, pp. 116–122, doi: 10.1109/ICCV.1999.791206 (1999).
- [2] J. Wu, Z. Tian, L. Sun, L. Estevez, R. Jafari, Real-time american sign language recognition using wrist-worn motion and surface emg sensors, in: *Wearable and Implantable Body Sensor Networks (BSN), 2015 IEEE 12th International Conference on*, IEEE, 2015, pp. 1–6, doi: 10.1109/BSN.2015.7299393 (2015).
- [3] S. Lahoti, S. Kayal, S. Kumbhare, I. Suradkar, V. Pawar, Android based american sign language recognition system with skin segmentation and svm, in: *2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, IEEE, 2018, pp. 1–6, doi: 10.1109/ICCCNT.2018.8493838 (2018).
- [4] S. Huang, C. Mao, J. Tao, Z. Ye, A novel chinese sign language recognition method based on keyframe-centered clips, *IEEE Signal Processing Letters* 25 (3) (2018) 442–446, doi: 10.1109/LSP.2018.2797228 (2018).
- [5] E.-J. Holden, G. Lee, R. Owens, Australian sign language recognition, *Machine Vision and Applications* 16 (5) (2005) 312, doi: <https://doi.org/10.1007/s00138-005-0003-1> (2005).
- [6] N. Tubaiz, T. Shanableh, K. Assaleh, Glove-based continuous arabic sign language recognition in user-dependent mode, *IEEE Transactions on Human-Machine Systems* 45 (4) (2015) 526–533, doi: 10.1109/THMS.2015.2406692 (2015).

- [7] T. Shanableh, K. Assaleh, M. Al-Rousan, Spatio-temporal feature-extraction techniques for isolated gesture recognition in arabic sign language, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 37 (3) (2007) 641–650, doi: [10.1109/TSMCB.2006.889630](https://doi.org/10.1109/TSMCB.2006.889630) (2007).
- [8] S. Hore, S. Chatterjee, V. Santhi, N. Dey, A. S. Ashour, V. E. Balas, F. Shi, Indian sign language recognition using optimized neural networks, in: *Information Technology and Intelligent Transportation Systems*, Springer, 2017, pp. 553–563, doi: https://doi.org/10.1007/978-3-319-38771-0_54 (2017).
- [9] F. López-Colino, J. Colás, Spanish sign language synthesis system, *Journal of Visual Languages & Computing* 23 (3) (2012) 121–136, doi: <https://doi.org/10.1016/j.jvlc.2012.01.003> (2012).
- [10] B. R. Barricelli, S. Valtolina, A visual language and interactive system for end-user development of internet of things ecosystems, *Journal of Visual Languages & Computing* 40 (2017) 1–19, doi: <https://doi.org/10.1016/j.jvlc.2017.01.004> (2017).
- [11] M. J. Cheok, Z. Omar, M. H. Jaward, A review of hand gesture and sign language recognition techniques, *International Journal of Machine Learning and Cybernetics* 10 (1) (2019) 131–153, doi: <https://doi.org/10.1007/s13042-017-0705-5> (2019).
- [12] B. S. Parton, Sign language recognition and translation: A multidisciplinary approach from the field of artificial intelligence, *Journal of deaf studies and deaf education* 11 (1) (2005) 94–101, doi: <https://doi.org/10.1093/deafed/enj003> (2005).
- [13] G. Fang, W. Gao, D. Zhao, Large-vocabulary continuous sign language recognition based on transition-movement models, *IEEE transactions on systems, man, and cybernetics-part a: systems and humans* 37 (1) (2007) 1–9, doi: [10.1109/TSMCA.2006.886347](https://doi.org/10.1109/TSMCA.2006.886347) (2007).
- [14] W. Gao, J. Ma, J. Wu, C. Wang, Sign language recognition based on hmm/ann/dp, *International journal of pattern recognition and artificial intelligence* 14 (05) (2000) 587–602, doi: [10.1142/S0218001400000386](https://doi.org/10.1142/S0218001400000386) (2000).
- [15] B. Khelil, H. Amiri, Hand gesture recognition using leap motion controller for recognition of arabic sign language, in: *3rd International conference ACECS'16*, 2016 (2016).
- [16] P. Kumar, H. Gauba, P. P. Roy, D. P. Dogra, Coupled hmm-based multi-sensor data fusion for sign language recognition, *Pattern Recognition Letters* 86 (2017) 1–8, doi: <https://doi.org/10.1016/j.patrec.2016.12.004> (2017).
- [17] F.-S. Chen, C.-M. Fu, C.-L. Huang, Hand gesture recognition using a real-time tracking method and hidden markov models, *Image and vision computing* 21 (8) (2003) 745–758, doi: [https://doi.org/10.1016/S0262-8856\(03\)00070-2](https://doi.org/10.1016/S0262-8856(03)00070-2) (2003).
- [18] M. Al-Rousan, K. Assaleh, A. Talaa, Video-based signer-independent arabic sign language recognition using hidden markov models, *Applied Soft Computing* 9 (3) (2009) 990–999, doi: <https://doi.org/>

- 10.1016/j.asoc.2009.01.002 (2009).
- [19] M. Maraqa, R. Abu-Zaiter, Recognition of arabic sign language (arsl) using recurrent neural networks, in: Applications of Digital Information and Web Technologies, 2008. ICADIWT 2008. First International Conference on the, IEEE, 2008, pp. 478–481, doi: 10.1109/ICADIWT.2008.4664396 (2008).
- [20] M. Mohandes, M. Deriche, Image based arabic sign language recognition, in: Signal Processing and Its Applications, 2005. Proceedings of the Eighth International Symposium on, Vol. 1, IEEE, 2005, pp. 86–89, doi: 10.1109/ISSPA.2005.1580202 (2005).
- [21] M. Mohandes, M. Deriche, U. Johar, S. Ilyas, A signer-independent arabic sign language recognition system using face detection, geometric features, and a hidden markov model, Computers & Electrical Engineering 38 (2) (2012) 422–433, doi: <https://doi.org/10.1016/j.compeleceng.2011.10.013> (2012).
- [22] K. M. Lim, A. W. Tan, S. C. Tan, A feature covariance matrix with serial particle filter for isolated sign language recognition, Expert Systems with Applications 54 (2016) 208–218, doi: <https://doi.org/10.1016/j.eswa.2016.01.047> (2016).
- [23] T. E. Starner, Visual recognition of american sign language using hidden markov models., Tech. rep., Massachusetts Inst Of Tech Cambridge Dept Of Brain And Cognitive Sciences (1995).
- [24] O. Koller, H. Ney, R. Bowden, Deep learning of mouth shapes for sign language, in: Proceedings of the IEEE International Conference on Computer Vision Workshops, 2015, pp. 85–91, doi: 10.1109/ICCVW.2015.69 (2015).
- [25] O. Koller, O. Zargaran, H. Ney, R. Bowden, Deep sign: Hybrid cnn-hmm for continuous sign language recognition, in: Proceedings of the British Machine Vision Conference 2016, 2016, doi: <http://epubs.surrey.ac.uk/812319/> (2016).
- [26] O. Koller, H. Ney, R. Bowden, Deep hand: How to train a cnn on 1 million hand images when your data is continuous and weakly labelled, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 3793–3802, doi: 10.1109/CVPR.2016.412 (2016).
- [27] O. Koller, S. Zargaran, H. Ney, R. Bowden, Deep sign: Enabling robust statistical continuous sign language recognition via hybrid cnn-hmms, International Journal of Computer Vision 126 (12) (2018) 1311–1325, doi: <https://doi.org/10.1007/s11263-018-1121-3> (2018).
- [28] O. Koller, C. Camgoz, H. Ney, R. Bowden, Weakly supervised learning with multi-stream cnn-lstm-hmms to discover sequential parallelism in sign language videos, IEEE transactions on pattern analysis and machine intelligence Doi: 10.1109/TPAMI.2019.2911077 (2019).
- [29] O. Koller, J. Forster, H. Ney, Continuous sign language recognition: Towards large vocabulary statistical recognition systems handling multiple signers, Computer Vision and Image Understanding 141 (2015) 108–125, doi: <https://doi.org/10.1016/j.cviu.2015.09.013> (2015).

- [30] U. von Agris, M. Knorr, K. Kraiss, The significance of facial features for automatic sign language recognition, in: 2008 8th IEEE International Conference on Automatic Face Gesture Recognition, 2008, pp. 1–6 (Sep. 2008). doi:10.1109/AFGR.2008.4813472.
- [31] A. Karami, B. Zanj, A. K. Sarkaleh, Persian sign language (psl) recognition using wavelet transform and neural networks, *Expert Systems with Applications* 38 (3) (2011) 2661–2667, doi: <https://doi.org/10.1016/j.eswa.2010.08.056> (2011).
- [32] A. Barkoky, N. M. Charkari, Static hand gesture recognition of persian sign numbers using thinning method, in: *Multimedia Technology (ICMT), 2011 International Conference on*, IEEE, 2011, pp. 6548–6551, doi: 10.1109/ICMT.2011.6002201 (2011).
- [33] M. Moghaddam, M. Nahvi, R. H. Pak, Static persian sign language recognition using kernel-based feature extraction, in: *Machine Vision and Image Processing (MVIP), 2011 7th Iranian*, IEEE, 2011, pp. 1–5, doi: 10.1109/IranianMVIP.2011.61215391 (2011).
- [34] A. A. Zare, S. H. Zahiri, Recognition of a real-time signer-independent static farsi sign language based on fourier coefficients amplitude, *International Journal of Machine Learning and Cybernetics* 9 (5) (2018) 727–741, doi: <https://doi.org/10.1007/s13042-016-0602-3> (2018).
- [35] S. G. Azar, H. Seyedarabi, University of tabriz persian sign language dataset (UoT-PSL), available at: https://asatid.tabrizu.ac.ir/Files/603_122fa2a0-989c-4124-b2f6-dfe2b5eb03ff.pdf (2016).
- [36] M. Brand, N. Oliver, A. Pentland, Coupled hidden markov models for complex action recognition, in: *Computer vision and pattern recognition, 1997. proceedings., 1997 ieee computer society conference on*, IEEE, 1997, pp. 994–999, doi: 10.1109/CVPR.1997.609450 (1997).
- [37] L. R. Rabiner, B.-H. Juang, An introduction to hidden markov models, *ieee assp magazine* 3 (1) (1986) 4–16 (1986).
- [38] L. R. Rabiner, A tutorial on hidden markov models and selected applications in speech recognition, *Proceedings of the IEEE* 77 (2) (1989) 257–286 (1989).
- [39] F. I. Bashir, A. A. Khokhar, D. Schonfeld, Object trajectory-based activity classification and recognition using hidden markov models, *IEEE transactions on Image Processing* 16 (7) (2007) 1912–1919, doi: 10.1109/TIP.2007.898960 (2007).