

Ethnicity Classification Based on Gait Using Multi-view Fusion

De Zhang^{1,2}, Yunhong Wang¹ and Bir Bhanu²
¹Intelligent Recognition and Image Processing Laboratory
Beihang University, Beijing 100191, China

zhangde@cse.buaa.edu.cn, yhwang@buaa.edu.cn

²Center for Research in Intelligent Systems
University of California, Riverside, CA 92521, USA

bhanu@cris.ucr.edu

Abstract

The determination of ethnicity of an individual, as a soft biometrics, can be very useful in a video-based surveillance system. Currently, face is commonly used to determine the ethnicity of a person. Up to now, gait has been used for individual recognition and gender classification but not for ethnicity determination. This paper focuses on the ethnicity determination based on fusion of multi-view gait. Gait Energy Image (GEI) is used to analyze the recognition power of gait for ethnicity. Feature fusion, score fusion and decision fusion from multiple views of gait are explored. For the feature fusion, GEI images and camera views are put together to render a third-order tensor $(x, y, view)$. A multilinear principal component analysis (MPCA) is used to extract features from tensor objects which integrate all views. For the score fusion, the similarity scores measured from single views are combined with a weighted SUM rule. For the decision fusion, ethnicity classification is realized on each individual view first. The classification results are then combined to make the final determination with a majority vote rule. A database of 36 walking people (East Asian and South American) was acquired from 7 different camera views. The experimental results show that ethnicity can be determined from human gait in video automatically. The classification rate is improved by fusing multiple camera views and a comparison among different fusion schemes shows that the MPCA based feature fusion performs the best.

1. Introduction

In human-centered video surveillance systems, soft biometric traits are fuzzy descriptions of the subjects under surveillance, such as gender, age, height, weight, eye color and ethnicity. These traits are important for tracking across

non-controlled cameras, surveillance monitoring, unconfident decision etc. Ethnicity is a fundamental characteristic of a person and lots of effort has been devoted to estimate ethnicity from face images over the past several years [4-6]. However, in many real-world video surveillance systems, it is hard to capture face information at a high enough resolution if the person is far away from the camera. It will lead to an unreliable identification of the ethnicity of a person. When a subject is far away from the camera, human gait can be detected and measured. This modality has the advantage at a distance when other biometric modalities might not be suitable due to the lack of resolution and non-cooperative subject [1]. So we propose a gait based ethnicity classification system in this paper. For optimal performance, an experimental environment in which several camera views are available is considered to capture as much gait information as possible.

A database which includes seven different views and two types of ethnicity is built for this study. Features are obtained using Gait Energy Image (GEI) from each view separately. We design three schemes for multi-view fusion. The first one is at feature level, in which we augment the dimension of sample space by taking camera view into account. Given the GEI of size m -by- n and the number of different camera views c , the original feature point in the sample space becomes a tensor object of size m -by- n -by- c . Then, in this third-order tensor space, a feature extraction method called multilinear principal component analysis (MPCA) [18] is implemented. The second scheme is at score level and the third one is at the decision level, when considering each single view as a classifier. Weighted sum rule is used to combine the matching scores and majority vote rule at decision level.

The paper is organized as follows. Section 2 introduces related work and the contribution of this paper. Section 3 presents the construction of GEI for all camera views and the fusion strategies. In Section 4, the database is described

in detail and experimental results are compared and discussed. Section 5 concludes the paper.

2. Related Work and Contribution

2.1. Related Work

Over the years there have been many successful studies about human identification and gender recognition based on gait. These studies provide us a valued guidance and reference to work on the issue of gait based ethnicity classification.

In [1] the authors present a review of databases and techniques for automatic gait recognition in a single view. The distinction between gaits of different genders is shown in [7-10]. SVM classifier is trained and tested on different gait features [7-9]. In [7] Lee et al. describes a representation which is comprised of parameters of moment features in 7 regions of a silhouette image. In [8] the gait signature is represented by a sequence of stick figure with 8 sticks and 6 joint angles. In [9] the authors use averaged gait image. Lawson et al. [10] propose a method of gait analysis utilizing the independent components of motion and demonstrate a high performance on gender classification by using the nearest neighbor classifier.

In recent years, multi-view gait recognition has been studied and the fusion of multi-view gait sequences generates improved results. Wang et al. [2] present a fusion scheme of multi-view gait sequences. They used CASIA database taken at 11 different views. Dempster-Shafer rule used at decision level produced a great improvement in comparison with single view based gait recognition. In [3] Huang et al. calculated the weight of each individual view by minimizing the probability of inaccurate classification and used these different weights to sum the distances between the test subject and the reference subject for each view. They used CMU MoBo database and chose five out of the six available viewing directions.

There are three possible levels of fusion as shown in [11]. One is to combine features at the feature extraction level, another is fusion at the matching score level and the last one is at the decision level. Table 1 shows a summary of fusion schemes of related work.

2.2. Contribution of This Paper

The contribution of this paper is three-fold. Firstly, we explore the problem of gait based ethnicity classification using a multi-view fusion. This is a new attempt in the field of soft biometrics analysis. Secondly, we propose a novel method to integrate gait information from multiple views at feature level. The GEI images as well as camera views generate a tensor sample space together. MPCA is used to extract feature from tensor data. Thirdly, we build a gait database specifically for ethnicity classification. Two types

of ethnicity and seven different camera views are available in this database. It is different from all existing gait databases. This database is used to compare different fusion techniques for gait based ethnicity classification.

Paper	Fusion level	Fusion method
Wang et al. [2]	Score level	<ul style="list-style-type: none"> • Sum rule • Product rule • D-S rule
Huang et al. [3]	Score level	Weighted sum rule
Zhou et al. [12]	Score level	<ul style="list-style-type: none"> • Sum rule • Product rule • Max rule
Zhou et al. [13]	Feature level	Feature concatenation and MDA
Shakhnarovich et al. [14]	Score level	Sum rule
Kale et al. [15]	Score level	<ul style="list-style-type: none"> • Hierarchical fusion • Sum rule • Product rule
Shan et al. [16]	Feature level	Canonical Correlation Analysis

Table 1. Fusion schemes used in related work.

3. Technical Approach

3.1. Feature Extraction

As described in [1], there are many approaches to extract gait features from 90° view by using the silhouette or designing a model. We need a representation that can be generalized to other view angles while characterizing gait effectively.

Gait Energy Image (GEI) is an effective representation of gait which reflects both static stance information of silhouettes and dynamic shape changing information over a gait cycle proposed by Han and Bhanu [20]. Yu et al. [19] use GEI algorithm [20] for 11 different views in CASIA gait database to evaluate the effect of view angle variation on gait recognition. In this work we also take GEI as an example to test the performance of recognizing ethnicity by gait. In order to construct GEI, a normalizing operation is required on the silhouettes extracted from original human walking videos with the technique of background subtraction. This includes scaling the foreground regions to the same height when keeping the ratio of its height to width and moving them to the center of silhouette images which have the same size. The next step in GEI construction is gait period detection. Let N_{gait} denote the number of frames included in one gait cycle. A simple strategy is proposed to estimate N_{gait} in [21]. But there exist some noise in the sig-



Figure 1. Examples of normalized and centered silhouette frames from different views for one walk. From the top row to bottom row, the view angles are 0° , 30° , 60° , 90° , 120° , 150° and 180° respectively. The rightmost image in each row is the corresponding gait energy image (GEI).

nal $N_f(t)$ which expresses the number of foreground pixels in the silhouette over time. Therefore, we use the autocorrelation technique to remove noise and estimate the length of a gait cycle accurately.

This method of detecting gait period works well for the gait sequences from the view of 90° . Because the cameras record human walking simultaneously in our database, the values of N_{gait} for different camera views are the same for one walk. The frame numbers decomposed from the original videos are used to synchronize the starting frame of a gait cycle for all the views in one walk. We extract one cycle of frames from the synchronized silhouettes for each camera view respectively and the grey-level gait energy image (GEI) can be calculated as:

$$G(x, y) = \frac{1}{N_{gait}} \sum_{t=1}^{N_{gait}} B_t(x, y) \quad (1)$$

where $B_t(x, y)$ is the silhouette image at time t in a gait period. Figure 1 shows some silhouette samples within a gait cycle from different views in our database and the rightmost images are the corresponding GEIs. As can be seen, GEI contains spatio-temporal information of a gait period in a compact way, no matter what the view angle is.

3.2. Fusion Schemes

Fusion can be done at three levels as mentioned in [11]. In this work, we explore the performance of using multiple views for gait based ethnicity classification on all these fusion levels.

3.2.1 Feature Level

A basic flowchart of the fusion scheme at feature level is shown in Figure 2(a). The method proposed here treats GEIs from seven different views as a data sample, which is a 3rd-order tensor. The spatial row space and column space of GEI as well as view space account for the 3 modes.

GEI is derived from a gait cycle as described above. We construct GEIs for all the different views as illustrated in Figure 1. Considering a walk recorded by cameras from different view angles simultaneously, we pick one gait cycle in this walk and obtain the corresponding GEIs for each view. A multi-view fusion can be realized at feature level when concatenating these two-dimensional GEIs in a three-dimensional space. Hence, we get the new integrated data samples in a 3rd-order tensor space shown in Figure 3.

In our database, there are 7 different view angles. As-

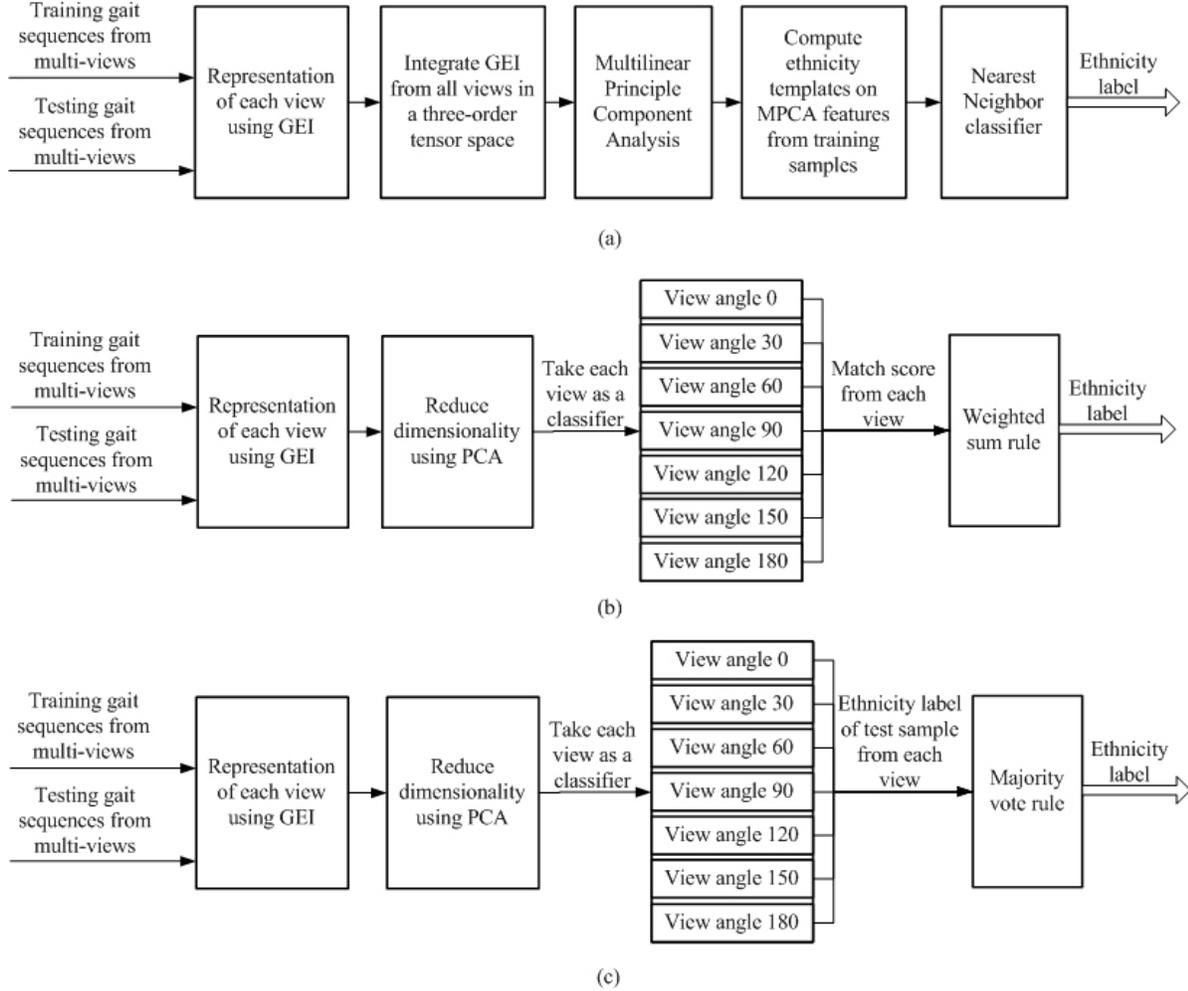


Figure 2. The basic frames of the fusion schemes used in this paper for comparison: (a) fusion at feature level, (b) fusion at score level, (c) fusion at decision level.

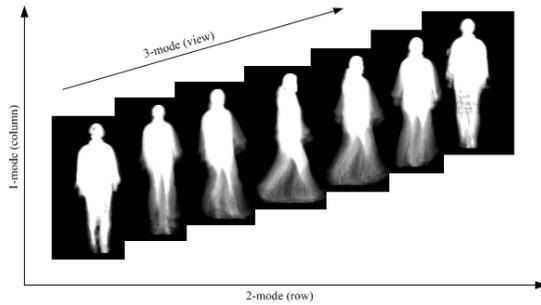


Figure 3. Illustration of GEIs from multiple views as a third-order tensor.

sume the size of GEI is m -by- n . Then, each walk is transformed into a tensor object of size m -by- n -by-7 using this fusion method at feature level. The whole data set to be an-

alyzed is a 4th-order tensor, with the addition of the sample space.

Lu et al. [18] introduce a multilinear principal component analysis (MPCA) framework for tensor object feature extraction. MPCA is a direct extension of PCA to the multilinear case. The input data samples are centered as in PCA, the projection is orthonormal and the projected feature is a tensor of the same order as the sample with reduced dimension. Given M training samples $\{\chi_1, \chi_2, \dots, \chi_M\}$ from N -order data set, compute its mode- N mean $\bar{\chi}^{(N)}$. The training data is centered by subtracting this mean tensor. Let χ^c denote the mode- N centered training samples. It can be decomposed using a higher-order SVD (HOSVD) as:

$$\chi^c = S \times_1 U^{(1)} \times_2 U^{(2)} \dots \times_2 U^{(2)} \quad (2)$$

As in SVD truncation for PCA, the HOSVD of χ^c is truncated by keeping in the first R_n columns for the basis matrix

$U^{(n)}$ in each mode n to produce $\tilde{U}^{(n)}$. The tensor projection using the collection of these $\prod_{n=1}^{N-1} R_n$ basis matrix is expressed as:

$$\tilde{p}_{-N} = \times_1 \tilde{U}^{(1)T} \times_2 \tilde{U}^{(2)T} \dots \times_{N-1} \tilde{U}^{(N-1)T} \quad (3)$$

These basis tensors are called eigentensors.

MPCA is an effective approach to deal with tensor data. In this paper, we apply it to reduce dimensionality of the tensor samples constructed from multiple views. The output of MPCA is the feature fused from multi-view gait.

3.2.2 Score Level

Fusion at score level is to combine the matching score from each system which indicates the proximity between the test feature vector and the template vector. The basic flowchart of the fusion scheme at score level is shown in Figure 2(b). We take each view as a classifier and compute the Euclidean distance between a test subject and the ethnicity template. This distance is the matching score.

Before combination, the scores from each view need to be normalized into a common domain. In this case, there are two matching scores for a given probe since our database includes only two types of ethnicity. So we transform these two values, S_1 and S_2 , in a simple way. The normalized scores are given by:

$$S_k^N = \frac{S_k}{S_1 + S_2}, k = 1, 2 \quad (4)$$

Then, a weighted SUM rule is used to combine the normalized scores from all views.

The weight should reflect the importance of each view. We resort to Fisher Linear Discriminant (FLD) to find appropriate value of weight assigned to an individual view. FLD maximizes the ratio of scatter between classes to the scatter within classes. This is done by maximizing $J(w)$ which is defined as:

$$J(w) = \frac{w^T S_B w}{w^T S_W w} \quad (5)$$

where S_B is scatter between classes matrix, S_W is scatter within classes matrix, and w is the projection vector. The maximum value of $J(w)$ indicates class separability for the given features. The training data from different camera views will result in different maximum values of $J(w)$. So, the weight of each view can be determined by the corresponding maximum value of $J(w)$. The total score s_t is calculated as:

$$J(w) = \sum_{i=1}^7 w_i s_i \quad (6)$$

where w_i is the weight of the i_{th} view and s_i is the corresponding score.

3.2.3 Decision Level

The basic flowchart of the fusion scheme at the score level is shown in Figure 2(c). Each classifier makes its own classification and votes for the final decision.

This strategy is motivated by the way humans make decisions, especially when there is a group of people involved in the decision process. Each classifier is in the position of a human expert with one vote. The resulting class is determined by the majority of votes. In our case, each view is taken as a single classifier and outputs an ethnicity label for a test subject. The majority vote rule is then used on these labels to determine the final ethnicity label.

4. Experiments

4.1. Database

The database is designed to include two types of ethnicity and multi-view gait sequences. Generally, the larger the distance between living areas is, the more difference exists between different ethnicities. So we found people from East Asia and South America to take part in the data collection. And we used eight cameras at seven different view angles to capture as much gait information as possible.

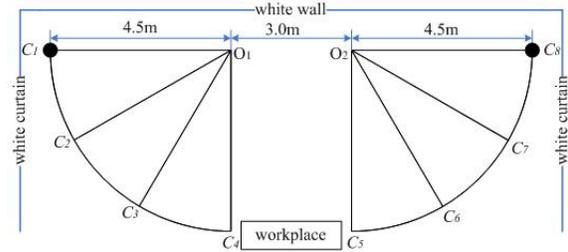


Figure 4. Cameras setup.

The gait video data were collected in an indoor environment. At the beginning of the collection process, the volunteers were asked to read, understand and sign an approved consent form. The form was a bilateral agreement for the usage rights of these data with the unique purpose of experiments. Then, they were asked to walk along a straight line five times between the two black solid points shown in Figure 4.

Eight cameras, from C_1 to C_8 as shown in Figure 4, were used in data collection. They were divided into two groups and in each group of four cameras formed a $1/4$ circle centered by O_1 and O_2 . The straight course between C_1 and C_8 was 12 meters approximately and it took about 9 seconds to walk such a distance at a normal speed. Every person walked along this straight line back and forth five times. We recorded their walking videos on two ways. As far as the view angles are concerned, we labeled C_1 with 180° , C_2 with 150° , C_3 with 120° , C_4 and C_5 with 90° , C_6 with 60° ,

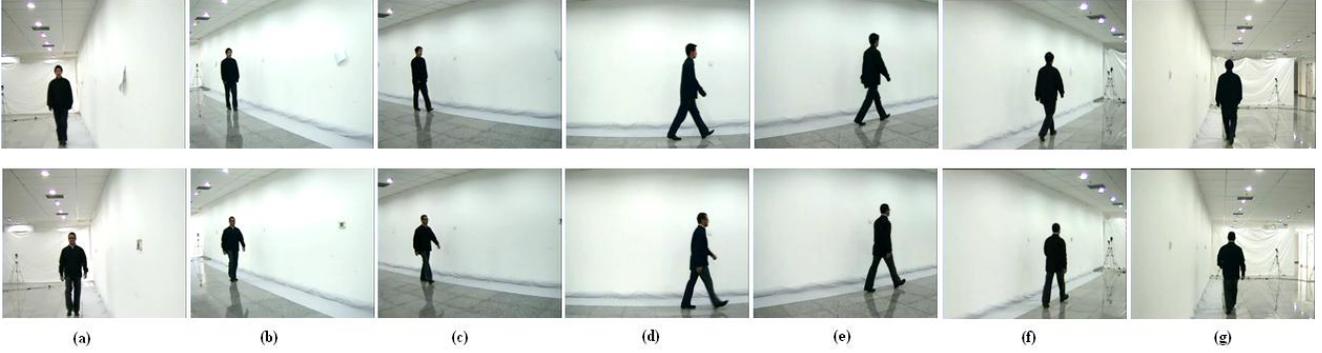


Figure 5. Frames at top row are from an East Asian subject. Frames at bottom row are from a South American subject. For each subject, these frames are synchronized in time. Columns (a) to (g) denote different view angles. (a) is from 0° , (b) is from 30° , (c) is from 60° , (d) is from 90° , (e) is from 120° , (f) is from 150° and (g) is from 180° .

C_7 with 30° , until C_8 with 0° . Note that we had C_4 and C_5 with the same view angle. It was for the purpose of increasing the distance between C_1 to C_8 so that the walking time could be longer in videos. Figure 5 shows the sample frames from all different camera views for two subjects of different ethnicities. The frames of video sequences are synchronized for these seven views.

The number of subjects in our database is 36 in total. There are 26 East Asians from China and Korea and 10 South Americans from Venezuela. They had a normal walk of five times at both left-to-right direction and the opposite direction recorded by eight cameras respectively during data collection. Hence, there are $2 \times 5 \times 8 = 80$ gait videos for each person in the database.

Since camera C_4 recorded the walking videos at the same view as camera C_5 , we left the data videos from C_4 for later experiments. Also, we only used data from the left-to-right walking direction. Therefore, each subject has the same five gait sequences for any view angle in our experiments.

4.2. Experimental Results

According to the fusion schemes discussed in Section 3, we performed three sets of experiment. To make the best use of data, we employed the leave-one-out cross-validation in all of these experiments. It involved using the overall five observations from one subject as the validation data, and the remaining observations as the training data. This was repeated such that each subject in the database was used once as the validation data. So, 36 iterations were required for 36 people in our database. We have 26 East Asians and 10 South Americans totally. The templates of the two types of ethnicity in training samples were built through computing the mean or median feature vector.

4.2.1 Experiment 1 - Feature Fusion

In this set of experiment, we implemented feature fusion as described in Section 3.2.1. The output of MPCA was transformed to a feature vector in which each factor was sorted in the order of decreasing eigenvalue. The recognition performance varied with the dimensionality of the MPCA feature vector. We selected the first k sorted features and tested their classifying performance with the nearest neighbor classifier. Figure 6 illustrates the correct classification rate versus MPCA dimensionality reduction.

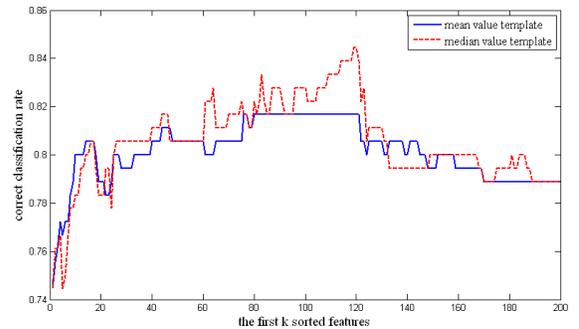


Figure 6. Different number of MPCA feature factors result in different classification rates.

For the median value based template, the classification rate achieves the peak at 120-dimension. For the mean value based template, it performs the best with 80-120 feature factors. In comparison, the maximum correct classification rate is higher when using median template.

4.2.2 Experiment 2 - Matching Score Fusion

In this set of experiment, we used the fusion scheme at matching score level as mentioned in Section 3.2.2. Given a test sequence of a single view, its matching scores related

to different ethnicity classes were calculated with Euclidean distance. Then, they were normalized to the same range [0, 1]. The normalized scores from different views were combined using the weighted SUM rule.

The weights were found by taking FLD analysis on PCA features of GEI. Of course, there was the issue of choose an appropriate value for the cumulative energy of eigenvectors. An empirical threshold was 90% in application of PCA. We used this value for the dimensionality reduction of GEI by PCA. Figure 7 shows the maximum values of $J(w)$ for all views. These values reflect the discriminating power of different views for ethnicity classification. So we took each of them as the weight for the corresponding view angle.

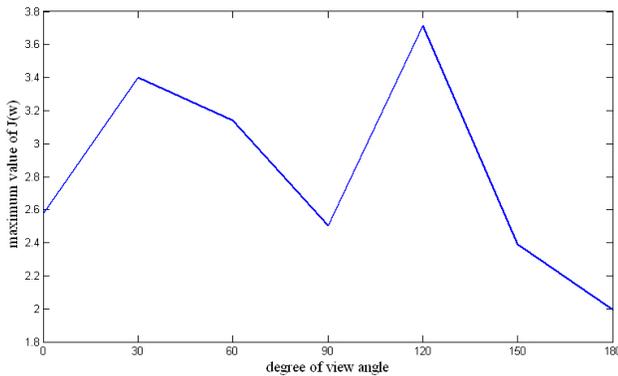


Figure 7. maximum $J(w)$ values.

The correct classification rates for mean template and median template are 79.4% and 80.5% respectively. For this fusion scheme, median template also gives a better result, the same case as fusion at feature level.

4.2.3 Experiment 3 - Decision Fusion

In this set of experiment, the fusion at decision level was explored. We used majority vote rule here. In this case, each view ran as an independent classifier and generated an ethnicity label. In one iteration in the leave-one-out cross-validation, one walk from a subject was chosen as the test data. Seven ethnicity classification labels related to this walk were produced from these seven view classifiers. Each label was a vote for the majority vote rule. The result of voting determined which ethnicity the person taking this walk belonged to. In our database, one person had five walks such that majority vote rule was used to make the final decision. Hence, the type of ethnicity of a subject used as the validation data in every repetition of the leave-one-out method was decided by applying the majority vote rule successively. The correct classification rates for mean template and median template are 78.1% and 78.5% respectively.

4.3. Discussion of Results

4.3.1 Classification Results for Each View

The correct classification rate for each single view is plotted in Figure 8. The median template produces the better results as it does in the case of multi-view fusion.

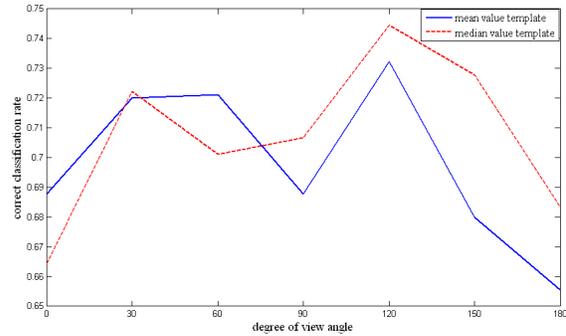


Figure 8. Classification results from each view.

Additionally, it is observed that the best classification rate among the single view based classification results is 74.4% at view of 120^0 . The view of 0^0 and 180^0 perform the worst. The similar case can be seen in Figure 7 which shows the different class separability of different single view. In the gait videos used in our experiments, the camera at 120^0 captures the oblique back of a subject. Our experimental results show this view is the most effective when classifying ethnicity.

4.3.2 Fusion Results

Table 2 shows the correct classification rates generated by fusion at different levels. As can be seen, median value based template gives a better result for any fusion scheme. The results show that the feature fusion performs the best. It is also indicated that GEI-based gait representation is able to classify ethnicity with a multi-view fusion, even using the simple Nearest Neighbor classifier.

Fusion level	Correct classification rate	
	Mean template	Median template
Feature level	81.7%	84.4%
Score level	79.4%	80.5%
Decision level	78.1%	78.5%

Table 2. Results of fusion at three levels.

The highest classification rate generated from our fusion schemes is 84.4%. It is much greater than the best one from the single 120^0 view. We achieved it by employing MPCA in the fusion at feature level. MPCA is able to produce a compact representation from the integrated multi-view gait information. It is a promising tool for gait analysis.

5. Conclusions

In this paper, gait based ethnicity classification is explored using multi-view gait fusion. A database consisting of two different ethnicities and seven different views is built for this task. A GEI is constructed from gait data acquired for each individual view and the GEIs are combined using different fusion schemes. Experimental results show that the performance of multi-view fusion is much better compared with a single view. The feature fusion scheme based on MPCA generates the best classification rate. This approach is a new proposal for multi-view gait analysis and it works well for the problem of ethnicity classification.

Acknowledgements

This work was supported by Natural Science Foundation of China (No. 60873158), 973 Program (No. 2010CB327902) and the opening funding of the State Key Laboratory of Virtual Reality Technology and Systems.

References

- [1] M. S. Nixon and J. N. Carter. Automatic recognition by gait. *In Proceedings of the IEEE*, vol. 94, issue 11, pp: 2013-2024, 2006.
- [2] Y. Wang, S. Yu, Y. Wang and T. Tan. Gait recognition based on fusion of multi-view gait sequences. *IAPR/IEEE International Conference on Biometrics*, LNCS 3832, pp: 605-611, 2005.
- [3] X. Huang and N. V. Boulgouris. Gait recognition using multiple views, *In Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, pp: 1705-1708, 2008.
- [4] S. Gutta, H. Wechsler and P. J. Phillips. Gender and ethnic classification of face images. *In Proc. of IEEE International Conference on Automatic Face and Gesture Recognition*, pp: 194-199, 1998.
- [5] X. Lu and A. K. Jain. Ethnicity identification from face images. *In Proc. of SPIE International Symposium on Defense and Security*, pp: 114-123, 2004.
- [6] Z. Yang and H. Ai. Demographic classification with local binary patterns. *IAPR/IEEE International Conference on Biometrics*, LNCS 4642, pp: 464-473, 2007.
- [7] L. Lee and W. E. L. Grimson. Gait analysis for recognition and classification. *In Proc. of IEEE International Conference on Automatic Face and Gesture Recognition*, pp: 148-155, 2002.
- [8] J. Yoo, D. Hwang and M. S. Nixon. Gender classification in human gait using support vector machine. *In Proc. of Advanced Concepts for Intelligent Vision Systems*, LNCS 3708, pp: 138-145, 2005.
- [9] X. Li, S. J. Maybank, S. Yan, D. Tao and D. Xu. Gait components and their application to gender recognition. *In IEEE Trans. on Systems, Man and Cybernetics, Part C: Applications and Reviews*, vol. 38, issue 2, pp: 145-155, 2008.
- [10] W. Lawson, Z. Duric and H. Wechsler. Gait analysis using independent components of image motion. *In Proc. of IEEE International Conference on Automatic Face and Gesture Recognition*, pp: 1-6, 2008.
- [11] A. Ross and A. K. Jain. Information fusion in biometrics. *Pattern Recognition Letters*, vol. 24, issue 13, pp: 2115-2125, 2003.
- [12] X. Zhou and B. Bhanu, Integrating face and gait for human recognition at a distance in video. *In IEEE Trans. on Systems, Man and Cybernetics, Part B: Cybernetics*, vol. 37, issue 5, pp: 1119-1137, 2007.
- [13] X. Zhou and B. Bhanu, Feature fusion of side face and gait for video-based human identification. *Pattern Recognition*, vol. 41, issue 3, pp: 778-795, 2008.
- [14] G. Shakhnarovich, L. Lee and T. Darrell. Integrated face and gait recognition from multiple views. *In Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, pp: 439-446, 2001.
- [15] A. Kale, A. K. RoyChowdhury and R. Chellappa. Fusion of gait and face for human identification. *In Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, vol. 5, pp: 901-904, 2004.
- [16] C. Shan, S. Gong and P. W. McOwan. Fusing gait and face cues for human gender recognition. *Neurocomputing*, vol. 71, issue 10-12, pp: 1931-1938, 2008.
- [17] H. Lu, K. N. Plataniotis and A. N. Venetsanopoulos. MPCA: multilinear principal component analysis of tensor objects. *In IEEE Trans. on Neural Networks*, vol. 19, issue 1, pp: 18-39, 2008.
- [18] S. Yu, D. Tan and T. Tan. A framework for evaluating the effect of view angle, clothing and carrying condition on gait recognition. *In Proc. of IEEE International Conference on Pattern Recognition*, vol. 4, pp: 441-444, 2006.
- [19] J. Han and B. Bhanu. Individual recognition using gait energy image. *In IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 28, issue 2, pp: 316-322, 2006.
- [20] S. Sarkar, P. J. Phillips, Z. Liu, I. R. Vega, P. Grother and K. W. Bowyer. The HumanID gait challenge problem: data sets, performance, and analysis. *In IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 27, issue 2, pp: 162-177, 2005.