# Recognising the Dynamics of Faces across Multiple Views

Yongmin Li, Shaogang Gong and Heather Liddell
Department of Computer Science
Queen Mary and Westfield College, University of London
London, E1 4NS, UK
{yongmin,sgg,heather}@dcs.qmw.ac.uk

### Abstract

We present an integrated framework for dynamic face detection and recognition, where head pose is estimated using Support Vector Regression, face detection is performed by Support Vector Classification, and recognition is carried out in a feature space constructed by Linear Discriminant Analysis. Unlike most traditional approaches to matching the patterns from static face images, we model the dynamics of human faces from video sequences in a consistent spatio-temporal context, i.e. recognition is accomplished by matching an object trajectory to a set of identity model trajectories in feature space. The model trajectories are synthesized from only a few views which sparsely cover the view sphere. Compared with the static face matching techniques, this approach is more robust and accurate under a coarse correspondence of face images, and has potential to visual interaction and advanced human behaviour recognition in real-world scenarios.

## 1 Introduction

The issue of face recognition has been extensively addressed over the past decade. Various models and approaches have been proposed aiming to solve the problem under different assumptions and constraints. Among them, the eigenface approach proposed in [24] uses Principal Component Analysis (PCA) to code face images and capture face features. This approach has then been extended to view-based and modular eigenspaces intended for recognising faces under varying views[14]. In [27], face recognition is performed by Elastic Graph matching based on a Gabor wavelet transform. An alternative approach uses similarity vectors to estimate head pose and recognise faces across views [10]. Active Shape Model (ASM) and Active Appearance Model (AAM) capturing both shape and shape-free grey-level appearance of face images have been successfully applied to face modelling and recognition [3, 4, 7]. Both ASM and AAM have been extended to nonlinear cases across views based on Kernel Principal Component Analysis (KPCA)[17, 19, 18]. These nonlinear models aimed at corresponding dynamic appearances of both shape and texture across views.

In most of the previous work, the basic methodology adopted for recognition is largely based matching static face image patterns in a given feature space. More recently, there has been some work on face recognition using video sequences [11, 8, 28, 20]. Nevertheless, the issue of recognising the dynamics of human faces under a spatio-temporal context remains largely unresolved.

In this paper, we describe an integrated approach to face detection and recognition with a focus on modelling and recognising the dynamics of faces in a spatio-temporal context. First, an overall framework is briefly introduced in Section 2. The methods and

process of head pose estimation and multi-view face detection are described in Section 3. Section 4 addresses the central issue of this paper: Recognising the dynamics of human faces. Conclusions are drawn in Section 5.

## 2 A Framework for Multi-view Face Detection, Tracking and Recognition

The aim is to model the dynamics of human faces undergoing rotations in depth, changes in scale, transformation in position, and variation to identities. We adopt as the input to our system live video sequences of captured but unsegmented faces that vary continuously in a temporal context. At the initializing phase, the system performs the follow tasks:

1. use motion estimation, skin colour detection, and background subtraction to bootstrap sub-images containing faces;
2. scan the sub-images with different scales;
3. for each scan, estimate the likely "pose" in tilt and yaw of the image patch;
4. choose a proper face detector from a set of multi-view face detectors according to the estimated "pose" to determine whether the pattern is a face. If the output of the face detector is above a preset threshold, then a face is detected, and the position, scale and pose of the detected face are used for recognition. Otherwise, the patch is rejected as a non-face pattern;
5. synthesize all detections to a single detection;
6. combine the detected face pattern, its pose estimate, and, if available, the history information from the previous frames to recognise the dynamics of the face.

After a face is successfully detected, it is not necessary to repeat this whole procedure again. For example, a Kalman filter can be used to track the position, scale, and pose of the detected face in the successive frame. When the tracking fails, the system can be restored by re-initializing as described above.



Figure 1: Multi-view face detection and recognition.

It is worth noticing that pose estimation is performed on multi-scale image patches before detection. If an image patch is detected as a face, then the tilt and yaw angles of the face is obtained, and the face detectors may give a positive output for the pattern. However, when the image patch is not exactly a face, the estimated pose is just a meaningless value. At the detection stage, one of the detectors will give a negative output to reject it as non-face. It seems that extra computation is exerted on non-face patterns at the stage of pose estimation. However, more saving in computation is achieved at the stage of multi-view face detection. The details will be described in Section 3. The proposed framework for dynamic face detection and recognition is illustrated in Figure 1.

# 3 Estimating Pose and Detecting Faces Across Views

The 2D appearance of faces across views changes significantly due to self-occlusion, change of illumination and nonlinear transformation. Obviously, this makes the problem of face recognition across views more challenging than that of a fixed view. However, if the view of a face image is known, then the recognition problem can be simplified to a great extent. Also, view information provides a useful cue for motion prediction, object tracking, and intention understanding, which are important in visual interaction and high-level behaviour recognition [21].

We use Support Vector Regression (SVR) [25, 5, 22] to perform head pose estimation. The training set in our experiments contained 2660 face images taken from 20 subjects, 133 views of each. The view information of each image is labelled with the angles in tilt and yaw. After normalization, the images were projected to a feature space constructed by Principal Component Analysis (PCA) [16, 15]. The PCA was also performed on the same set of images. Then the projected patterns were fed into a decomposition SVR algorithm for training.

Two SVR based pose estimators, one for yaw and the other for tilt, were trained to estimate the head poses. In our experiments, the proportion of SVs is only about 5% of the training examples. More details of how to construct the pose estimators are given in [13].

Face detection can be defined as a classification problem of separating face patterns from non-face patterns. There are basically three methods to perform multi-view face detection. A straightforward way is to build a single detector dealing with all views of faces. However, due to nonlinear variations of faces between different views which makes the distribution of patterns much irregular, this approach leads to poor performance in most cases. The second approach is to build several detectors, each of them corresponding to a specific view. When detecting, all the detectors are employed, and if one or more of them gives positive output, then a face is considered to be detected. From our experiments, this approach performs better, but, as expected, the computation is more expensive, as each of the multi-view face detectors is calculated on a given pattern. Also, how to synthesize the final detection from the outputs of all the detectors is a nontrivial problem. We adopted a third approach which employs head pose information explicitly.

We constructed four face detectors using Support Vector Classification (SVC) [2, 25]. When training each of the multi-view face detectors, the face images corresponding to the specific views are selected from the same database for pose estimation as positive examples (faces). Negative examples (non-faces) are collected by a boot-strapping method [23] from a set of scenic pictures.

The size of the example images is 20x20 in pixels. After normalization, the images are projected into the PCA space as described in the previous section. It is worth noticing that the results of PCA were trained only on positive images (faces) since we are only interested in detecting faces.

Provided the pose of an image patch is known, the computation of face detection can be greatly reduced by only choosing the appropriate face detector for the given pose [13].

# 4 Recognising the Dynamics of Faces

## 4.1 Linear Discriminant Analysis

Linear Discriminant Analysis (LDA) seeks to find a linear transformation from input space to a feature space (LDA space) by maximizing the between-class variance and

minimizing the within-class variance at the same time. Computationally, LDA is similar to eigen-decomposition [9]. Since it is quite effective to select discriminant features for two or more object classes, LDA has been widely adopted in many pattern recognition applications. Zhao et al. [29] used LDA as a representation for frontal-view face recognition. Edwards et al. [6] adopted LDA to select *Discriminant Parameters* based on Active Appearance Models.

In this paper, we adopt a similar approach to [29] for the representation of face images. However we extend the problem to *multi-view* face recognition. Pose change is between -90° -90° in yaw and -30° -30° in tilt in all experiments of this work. Intuitively, the image appearance of different people at the same view is more similar than that of the same person at different views. This makes LDA less effective when applied to multi-views. This is illustrated in Figure 2(a), where the distribution of face patterns from 3 subjects, 133 views of each, is showed in the LDA feature space spanned by the first 3 significant LDA features. One can observe the fact that the variance from different subjects is not significantly isolated from the variance from different views, which suggests that LDA alone is insufficient for multi-view face recognition.



Figure 2: Distribution of multi-view face patterns in the LDA space spanned by the first 3 LDA features. (a) Patterns from 3 subjects, 133 views of each. (b) Selected patterns by pose from the left figure, where yaw changes from -90° to 90° and tilt is 0° for all patterns.

However, if we consider the pose information, the problem becomes different. Figure 2(b) shows the selected patterns from Figure 2(a) where tilt is 0° and yaw changes from -90° to 90° . We form the trajectories of the patterns from the same subjects. One notices that, even under a very low dimensional feature space (only the first 3 LDA features used), the identities of different subjects captured by their trajectories are separable, therefore can be discriminated.

## 4.2   Synthesizing Virtual Views for Multi-view Face Recognition

If all views of subjects are available, the problem of face recognition can be simply matching the face patterns on a specific view. However, this assumption may be too strict in many situations. That leads the following question: Given a few views of a subject, can we still perform recognition on unknown views? One solution is to synthesize virtual views in some feature space based on face patterns available, and then to match the unseen face pattern to the synthesized patterns. Vetter and Poggio [26] used Linear Object Classes for this purpose. In their work, the synthesis of virtual views was based on a dense correspondence between 2D face images of different people at the same view. Alternatively similarity vectors to a set of prototypes were adopted to construct the feature space, where face recognition on novel views was performed by linear interpolation of similarity vectors of patterns on available views [10].

In this paper, we present an approach to synthesize virtual views in LDA space. The problem can be formed as follows: Given a sparse set of face patterns $P_1, P_2, \cdots, P_m$ of an object on prototype views $v_1, v_2, \cdots, v_m$, one approximates the distributions of this object in feature space $P(v)$ by the patterns available. Figure 3 shows real distributions of face patterns from one subject with respect to different views and the synthesized distribution from 15 out of 133 views.



Figure 3: Comparison between original face patterns and synthesized patterns from 15 views. The three axes are tilt, yaw, and the first 3 LDA features respectively.

If more prototype views are available, a high order approximation of $P(v)$ is preferable for high accuracy. However, if only a few prototype views are available, one can still perform the synthesis by simple linear techniques such as bilinear transform which approximates a novel view by the four nearest prototype views, or by a triangle transform which uses the 3 nearest prototype views.

In our system, after face detection, the position, scale, and pose both in tilt and yaw of the detected face are obtained. Then for each subject to be recognised, a synthesized pattern in the detected view is generated in the feature space from the prototype views of the corresponding subject. Finally recognition is performed by matching the detected pattern with all the synthesized patterns.

| n | tilt | yaw | rate | known | novel |
|---|------|-----|------|-------|-------|
| 4 | 80,100 | 40,140 | 68.54% | 74.15% | 62.92% |
| 6 | 80,100 | 40,90,140 | 82.92% | 87.08% | 78.77% |
| 9 | 70,90,110 | 30,90,150 | 90.00% | 94.62% | 85.38% |
| 15 | 70,90,110 | 30,60,90,120,150 | 94.31% | 96.62% | 92.00% |

Table 1: Recognition results of using synthesized patterns in LDA feature space. n: number of prototype patterns; tilt/yaw: view angles of prototype patterns; rate: overall recognition accuracy; known/novel: accuracy on known/novel subjects respectively.

We trained LDA on a set of face images from 10 subjects, 133 views of each. The test set included face images from 20 subjects, where 10 subjects did not appear in the training set. Table 1 lists the results of recognition using different selections of prototype views. In this experiment, we did not perform pose estimation on the test images, instead the ground-truth pose information was used for synthesizing virtual views. The results are encouraging as they depict that:

1. given accurate pose estimation and face detection, a small number of prototype views are sufficient to synthesize novel views for recognition;

2. the LDA representation obtained from only 10 subjects generalizes well on novel subjects since the recognition accuracy on novel subjects is only slightly lower than that on known subjects.

## 4.3   Recognition Sensitivity to Image Alignment and Pose

It is crucial to point out that *no explicit alignment is performed in our approach*. In fact, face detection takes the role of coarse correspondence based on the appearance of image patterns. Compared with other methods such as optical flow or Active Shape Models where the correspondence between images is aimed to be well established either densely or sparsely, this method is only based on registering the holistic patterns approximately. Also, the synthesis of virtual views is based on pose estimation where the average error in our system is around $10°$ in both tilt and yaw. This error is another factor which influences the accuracy of recognition.



|  (a)  |  (b)  |

Figure 4: Sensitivity analysis with respect to position shift of detected faces. (a) Recognition rate with respect to the shift along horizontal/vertical direction of images. (b) Recognition rate of subjects being recognised in the top $n$ matches. Only 6 curves are plotted here. The dotted curve is from the results of the accurate position. The corresponding shifts for the other 5 curves are -2 in $x$, and -2 to 2 in $y$. In this experiment, 15 prototype views as shown in Table 1 were used, and the face image size is 20x20.

We performed sensitivity analysis based on two factors: Sensitivity to the image shift off the centre, and sensitivity to pose estimate. Figure 4 shows the recognition rates of the top match (the nearest distance) when shifting the ground-truth positions of faces around centres and the recognition rates obtained from the top $n$ matches (the first $n$ nearest distances). Figure 5 gives the recognition results when all the view information is estimated by our pose estimators with comparison to the results on ground-truth pose information.

In Figure 4 and Figure 5, we noted: First, the error both in position alignment and view estimation have significant effects on the top recognition match *alone*. Second, most of the patterns that were not correctly recognised by the top match were correctly recognised within the top 3-5 matches.

If we improve the accuracy in both detection and pose estimation, the performance would be improved to some extent. However, this assumption may be too strict in many situations. An alternative solution is to incorporate explicit shape models into the system for an accurate correspondence, but usually this approach is expensive and can also be too restrictive in real-world dynamic scenes.

If recognition is performed on a continuously moving face by accumulated identification evidence rather than snapshots of the face independently, the requirement for alignment and pose estimation can be relaxed. Based on this, we present an approach to address the problem by recognising face dynamics.

Figure 5: Sensitivity analysis to the error in pose estimation. The two figures show the results when 4 and 6 prototype views were used for recognition respectively (see Table 1). The horizontal axis stands for the number of top matches, and vertical axis for recognition rate. Solid curves are obtained from estimated pose, and dotted curves from ground-truth pose.

## 4.4 Recognising Face Dynamics by Matching Trajectories in Feature Space

In many real-world applications, a sequence of images containing the subjects to be recognised are necessarily acquired. With a sequence which records a face varying continuously over time and across views, not only more information can be obtained, but also the dynamics of faces can be captured [12, 1].

Yamaguchi et al. [28] presented a method for face recognition from sequences by building a subspace for the detected faces on the given sequence and then matching the subspace with prototype subspaces. Gong et al. [11] introduced an approach that uses Partially Recurrent Neural Networks (PRNNs) to recognise temporal signatures of faces.

We present here an approach to recognising the dynamics of human face by matching trajectories in LDA feature space. For a given sequence containing faces to be recognised, after head pose estimation and face detection, one can obtain a trajectory by projecting the face patterns into LDA space. On the other hand, according to the pose information of the face patterns, it is easy to build the identity model trajectory for each subject using the known prototype patterns. Therefore, the recognition problem can be solved by matching the object trajectory to a set of identity model trajectories.

At frame $t$ of a sequence, we define the distance between the object trajectory and a identity model trajectory $m$ as:

$$d_m = \sum_{i=1}^{t} w_i d_{mi} \qquad (1)$$

where $d_{mi}$ is the distance in feature space between model point and object point at frame $i$, and $w_i$ is the weight on this distance. Considerations on determining $w_i$ include:

1. confidence of face detection;
2. variation from the previous frame;
3. view of the face pattern. For example, profile face patterns are weighted lower than those at 3/4 views since they carry less discriminating information.

Finally, result of recognition can be given as:

$$id = argmax_{m=1}^{M} d_m \qquad (2)$$

Figure 6 illustrates an example of recognising the dynamics of moving faces. Recognition through matching static face patterns in individual frame can lead to wrong results.

However, if recognition is performed by matching the accumulated trajectories, a more robust and accurate recognition can be achieved. In particular, it is important to observe in Figure 6(c,d) that the true identity of a facial appearance is often recognised among the best few matches although it may not be consistently the best match in every frame over time. In other word, the accumulation of positive identity information will overwhelm any misidentification over time if recognition is performed on accumulated evidence.



(a)



(b)           (c)           (d)

Figure 6: Recognising face dynamics by matching trajectories in LDA feature space. (a) Successive frames from a test sequence where detected faces are marked with white boxes. (b) Object and identity model trajectories. The object trajectory is showed by solid curve with frames labelled by small circles. The others are model trajectories. For clarify, only 3 out of 20 model trajectories are illustrated here. (c) Distance measured independently in each frame between object pattern and model patterns. (d) Distance between object and different model trajectories over time. Results on the first 5 subjects are illustrated here, where the solid line is from the true subject.

# 5 Conclusions

From the viewpoint of visual interaction and human-computer interaction, the problem of face recognition involves more than matching static images. At a low-level, the face dynamics can be accommodated in a consistent spatio-temporal context where the underlying variations with respect to changes in identity, view, scale, position, illumination, and occlusion are integrated together. At a higher level, more sophisticated behaviour models, including individual-dependent and individual-independent models, may supervise and co-operate with all the low level modules.

In this work, we described an integrated framework for head pose estimation, multi-view face detection and the recognition of face dynamics. Key issues of the work and main differences from the previous work [6, 29, 11, 20, 28] include:

1. It deals with faces undergoing large rotation in depth, i.e. -90° -90° in yaw and -30° -30° in tilt. Head pose estimation by Support Vector Regression, and the pose information is explicitly used to simplify and guide the processes of face detection, recognition and face dynamics analysis.

2. The problem of multi-view face detection is decomposed into a set of sub-problems on different small ranges of views based on the pose information. Support Vector Classification is employed to solve those sub-problems.

3. *Coarse registration* of face patterns is adopted to avoid extensive computation in the process of correspondence.

4. No constraint is imposed in this work that all views of a subject should be available. In fact, by synthesizing virtual views, only a few views are sufficient for capturing the whole view sphere. Furthermore, those prototype views need not to be fixed, a randomly sampled sequence clip can be conveniently used for modelling.

5. Face Recognition is performed by modelling face dynamics in LDA feature space other than matching static images.

We believe that exploiting the dynamics of human faces is the most promising issue in face recognition. It is true that measuring the distance between the object and identity model trajectories is only a simple implementation of the approach. More elaborated methods, such as matching the temporal variation of the trajectories using higher order parameters other than only positions in feature space, can be more effective. In this paper, we highlighted the nature of the problem and showed the potential of modelling face dynamics as an effective means to face recognition. Nevertheless, we only touched the surface of the problem by modelling the dynamics across views and between subjects in an LDA feature space. Substantial future work using sophisticated temporal models to capture face dynamics with respect to expression, movement and illumination changes is to be conducted.

# References

[1] V. Bruce, A. Burton, and P. Hancock. Comparisons between human and computer recognition of faces. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 408–413, Nara, Japan, 1998.

[2] C. J. C. Burges. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2(2):1–47, 1998.

[3] T. Cootes, G. Edwards, and C. Taylor. Active appearance models. In *European Conference on Computer Vision*, volume 2, pages 484–498, Freiburg, Germany, 1998.

[4] T. Cootes and C. Taylor. A mixture model for representing shape variation. *Image and Vision Computing*, 17:567–573, 1999.

[5] H. Drucker, C. J. C. Burges, L. Kaufman, A. Smola, and V. Vapnik. Support vector regression machines. In M. Mozer, M. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems 9*. MIT Press, Cambridge, MA, 1997.

[6] G. Edwards, A. Lanitis, C. Taylor, and T. Cootes. Statistical models of face images - improving specificity. In *British Machine Vision Conference*, volume 2, pages 765–774, Edinburgh, Scotland, 1996.

[7] G. Edwards, C. Taylor, and T. Cootes. Interpreting face images using active appearance models. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 300–305, Nara, Japan, 1998.

[8] G. Edwards, C. Taylor, and T. Cootes. Learning to identify and track faces in sequences. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 260–267, Nara, Japan, 1998.

[9] K. Fukunaga. *Introduction to statistical pattern recognitiion*. Academic Press, 1972.

[10] S. Gong, E.-J. Ong, and S. McKenna. Learning to associate faces across views in vector space of similarities to prototypes. In *British Machine Vision Conference*, pages 54–64, Southampton, England, 1998.

[11] S. Gong, A. Psarrou, I. Katsouli, and P. Palavouzis. Tracking and recognition of face sequences. In *European Workshop on Combined Real and Synthetic Image Processing for Broadcast and Video Production*, pages 96–112, Hamburg, Germany, 1994.

[12] S. Gong and A. P. S. McKenna. *Dynamic Vision: From Images to Face Recognition*. World Scientific Publishing and Imperial College Press, April 2000.

[13] Y. Li, S. Gong, and H. Liddell. Support vector regression and classification based multi-view face detection and recognition. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 300–305, Grenoble, France, 2000.

[14] B. Moghaddam and A. Pentland. Face recognition using view-based and modular eigenspaces. In *Automatic Systems for the Identification and Inspection of Humans, SPIE*, volume 2277, 1994.

[15] H. Murase and S. K. Nayar. Visual learning and recognition of 3-d objects from appearance. *International Journal of Computer Vision*, 14:5–24, 1995.

[16] A. Pentland, B. Moghaddam, and T. Starner. View-based and modular eigenspaces for face recognition. In *IEEE CVPR*, pages 84–91, Seatle, 1994.

[17] S. Romdhani, S. Gong, and A. Psarrow. A multi-view nonlinear active shape model using kernel pca. In *British Machine Vision Conference*, pages 483–492, Nottingham, UK, 1999.

[18] S. Romdhani, S. Gong, and A. Psarrow. On utilising template and feature-based correspondence in multi-view appearance models. In *European Conference on Computer Vision*, Dublin, Ireland, June 2000.

[19] S. Romdhani, S. Gong, and A. Psarrow. A generic face appearance model of shape and texture under very large pose variations from profile to profile views. In *International Conference on Pattern Recognition*, Barcelona, Spain, September 2000.

[20] S. Satoh. Comparative evaluation of face sequence matching for content-based video access. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 163–168, Grenoble, France, 2000.

[21] J. Sherrah, S. Gong, J. Howell, and H. Buxton. Interpretation of group behaviour in visually mediated interaction. In *International Conference on Pattern Recognition*, Barcelona, Spain, September 2000.

[22] A. Smola, B. Scholkopf, and K.-R. Muller. General cost functions for support vector regression. In T. Downs, M. Frean, and M. Gallagher, editors, *Proc. of the Ninth Australian Conf. on Neural Networks*, pages 79–83, Brisbane, Australia, 1998.

[23] K. Sung and T. Poggio. Example-based learning for view-based human face detection. Technical report, Massachusetts Institute of Technology, 1994. A.I. MEMO 1521.

[24] M. Turk and A. Pentland. Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3(1):71–86, 1991.

[25] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer Verlag, New York, 1995.

[26] T. Vetter and T. Poggio. Linear object classes and image synthesis from a single example image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):733–742, 1997.

[27] L. Wiskott, J. Fellous, N. Kruger, and C. Malsburg. Face recognition by elastic bunch graph matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):775–779, 1997.

[28] O. Yamaguchi, K. Fukui, and K. Maeda. Face recognition using temporal image sequence. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 318–323, Nara, Japan, 1998.

[29] W. Zhao, A. Krishnaswamy, R. Chellappa, D. Swets, and J. Weng. Discriminant analysis of principal components for face recognition. In Wechsler, Philips, Bruce, Fogelman-Soulie, and Huang, editors, *Face Recognition: From Theory to Applications*, pages 73–85. Springer-Verlag, 1998.