

A Model Based Approach for Expressions Invariant Face Recognition

Zahid Riaz, Christoph Mayer, Matthias Wimmer,
Michael Beetz, and Bernd Radig

Department of Informatics, Technische Universität München
85748 Garching, Germany
{riaz,mayerc,matthias.wimmer,beetz,radig}@in.tum.de

Abstract. This paper describes an idea of recognizing the human face in the presence of strong facial expressions using model based approach. The features extracted for the face image sequences can be efficiently used for face recognition. The approach follows in 1) modeling an active appearance model (AAM) parameters for the face image, 2) using optical flow based temporal features for facial expression variations estimation, 3) and finally applying classifier for face recognition. The novelty lies not only in generation of appearance models which is obtained by fitting active shape model (ASM) to the face image using objective functions but also using a feature vector which is the combination of shape, texture and temporal parameters that is robust against facial expression variations. Experiments have been performed on Cohn-Kanade facial expression database using 62 subjects of the database with image sequences consisting of more than 4000 images. This achieved successful face recognition rate up to 91.17% using binary decision tree (BDT), 98.6% using Bayesian Networks (BN) with 10-fold cross validation in the presence of six different facial expressions.

Keywords: Active Appearance Models, Face Recognition, Facial Expressions Recognition, Binary Decision Trees, Bayesian Classifier.

1 Introduction

Since last three decades of face recognition technology there exists many commercially available systems to identify human faces, however face recognition is still an outstanding challenge against different kinds of variations like facial expressions, poses, non-uniform light illuminations, occlusions and aging effects. Meanwhile this technology has extended its role to Human-Computer-Interaction (HCI) and Human-Robot-Interaction (HRI). Person identity is one of the key tasks while interacting with the robots, exploiting the oblivious system security and authentication of the human interacting with the system. In such kinds of scenarios the acquired face images could contain various types of facial expressions along with above mentioned effects. This problem has been addressed in this paper resulting in a face recognition system which is robust against facial expressions. In [1, 2] authors give comprehensive survey of face recognition and some commercially available face recognition software.

As cameras are widely used and mounted on computer screens, embedded in mobiles and installed into everyday living and working environments, they have become valuable tools for human system interaction. A particular important aspect of this interaction is detection and recognition of faces and interpretation of facial expressions. These capabilities are deeply rooted in the human visual system and a crucial building block for social interaction. Consequently, these capabilities are an important step towards the acceptance of many technical systems. Although faces are the most important and natural way for human-human interaction but some outstanding challenges like uniqueness and accuracy made its market value a bit less than the other biometrics in year 2003. Later in 2006 Face recognition technology again raised up to 19% of the biometric market [23].

This publication focuses on one of the aspects of natural human-computer interfaces: our goal is to build a real-time system for face recognition that could robustly run in real-world environments. We develop it using model-based image interpretation techniques, which have proven its great potential to fulfill current and future requests on real-world image understanding. Our approach comprises methods that robustly localize facial features, seamlessly track them through image sequences, and finally infer the face identity.

2 Related Work

Traditional recognition systems have the abilities to recognize the human using various techniques like feature based recognition, face geometry based recognition, classifier design and model based methods. Linear subspace methods like Principal Components Analysis (PCA) was firstly used by Sirvovich and Kirby [3], which were latterly adopted by M. Turk and A. Pentland introducing the famous idea of eigenfaces [4,5]. This paper focuses on the modeling of human face using a two dimensional approach of shape, texture and temporal information and then utilizing this model for recognition purposes. This type of model using shape and texture parameters is called Active Appearance Models (AAMs), introduced by Edwards and Cootes [6]. In [7] authors have used weighted distance classifier called Mahalanobis distance measure. However, Edwards et al [8] isolated the sources of variation by maximizing the interclass variations using Linear Discriminant Analysis (LDA), the technique which was holistically used for Fisherfaces representation [9] which is similar to the eigenface approach resulting in out performance of previous approach. In [10] authors have utilized shape and temporal features collectively to form a feature vector for facial expressions recognition. Lanitis et al [11] used separate information for shape and gray level texture. These models utilize the shape information based on a point distribution of various landmarks points marked on the face image. In our approach a predefined shape model consisting of 134 points in two dimensional space, as was used by Wimmer et al [12], has been utilized.

3 Our Approach

This paper utilizes shape model fitting, extracting shape information, texture mapping for texture parameterization and estimating optical flow-based parameters for feature

vector formation. The feature vector consists of the shape, texture and temporal variations, sufficient for considering local variations in shapes. All the subjects in the database are labeled for identification. A fitted face model, on the training images is then used for defining the reference shape in our experiments. This reference shape is calculated by finding the mean shape of the all shapes in the database.

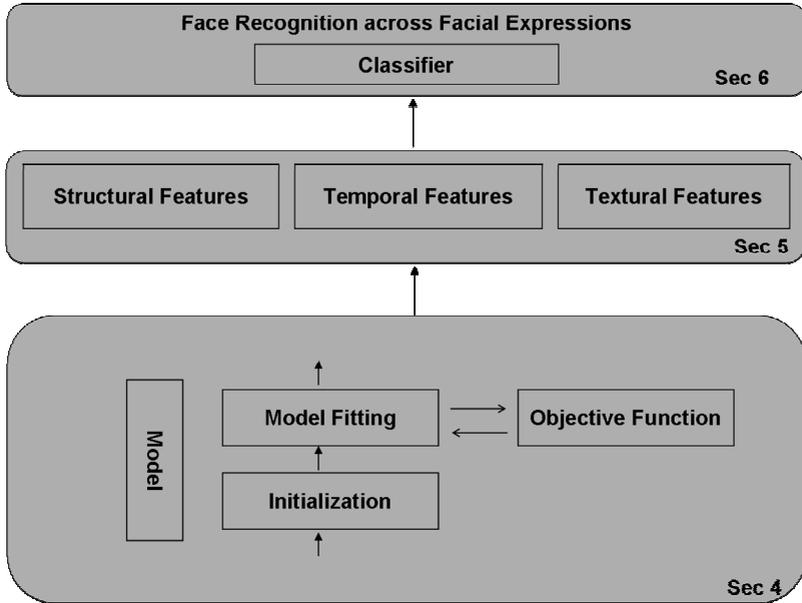


Fig. 1. Model-based techniques split the challenge of image interpretation into computationally independent modules. The lower right corners refer to the sections with detailed explanation.

In this paper an explicit 2D appearance model is used to develop a baseline for feature extraction. A point distribution model (PDM) is used as an active shape model. This model comprises of 134 points that prominently defines the location of local face features like eyes, nose and lips in 2D space. Face is localized in the image. An objective function is learned for fitting this model to the faces. After fitting the model to the example face image, texture information is extracted from the example image on a reference shape which is the mean shape of all the shapes of database. Image texture is extracted using planar subdivisions of the reference and the example shapes. Texture warping between the subdivisions is performed using affine transformation. This image texture is now normalized both in the sense of shape and varying illuminations effects, making the image robust for shape and illumination. Principal Component Analysis (PCA) is used to obtain the texture and shape parameters of the example image. This approach is similar to extracting Active Appearance Model (AAM) parameters. In addition to AAM parameters, temporal features of the facial changes are also calculated. Local motion of the feature points is observed using optical flow. We use reduced descriptors by trading off between accuracy and runtime performance. These features are then used for classifiers for

face recognition. A detailed process flow of our approach is shown in Figure 1. Our approach achieves real-time performance and provides robustness against facial expressions for real-world applicability. This computer vision task comprises of various phases for which it exploits model-based techniques that accurately localize facial features, seamlessly track them through image sequences, and finally infer facial expressions visible. We specifically adapt state-of-the-art techniques to each of these challenging phases.

The remainder of this paper is divided in four sections. In section 4 model based image interpretation is described along with the model fitting process. Sections 5 discuss about the model based feature extraction technique comprising shape and appearance along with the temporal features. Section 6 deals with feature classification. Finally, the experimentation is performed in section 7 along with conclusions.

4 Model-Based Image Interpretation and Fitting

Model-based techniques consist of four components: the model, the initialization algorithm, the objective function, and the fitting algorithm. Our approach makes use of a statistics-based deformable model, as introduced by Cootes et al. [13]. The model contains a parameter vector p that represents its possible configurations, such as position, orientation, scaling, and deformation. Models are mapped onto the surface of an image via a set of feature points, a contour, a textured region, etc. Referring to [6], deformable models are highly suitable for analyzing human faces with all their individual variations. Its parameters $p = (tx, ty, s, \theta, b)^T$ comprise the translation, scaling factor, rotation, and a vector of deformation parameters $b = (b_{s,1}, \dots, b_{s,m})^T$. The latter component describes the configuration of the face, such as the opening of the mouth, roundness of the eyes, raising of the eyebrows.

The initialization algorithm automatically starts the interpretation process by roughly localizing the object to interpret. It computes an initial estimate of the model parameters that needs to be further refined by the subsequent fitting algorithm. Our system integrates the approach of Viola and Jones, which is able to detect the affine transformation parameters (tx , ty , s , and θ) of frontal faces.

In order to obtain higher accuracy, we apply a second iteration of the Viola and Jones object detector to the previously determined image region of interest i.e. the face. This extension allows to roughly estimate the deformation parameters b_s , as well, because we learn the algorithm to localize facial components, such as eyes and mouth. In the case of the eyes, our positive training examples contain the images of eyes, whereas the negative examples consist of image patches in the vicinity of the eyes, such as the cheek, the nose, or the brows. Note that the resulting eye detector is not able to robustly localize the eyes in a complex image, because it usually contains a lot of information that was not part of the training data. However, it is highly appropriate to determine the location of the eyes within a pure face image or within the face region of a complex image.

The objective function $f(I, p)$ yields a comparable value that specifies how accurately a parameterized model p describes an image I . It is also known as the likelihood, similarity, energy, cost, goodness, or quality function. Without losing generality, we consider lower values to denote a better model fit. Traditionally, objective functions are manually specified by first selecting a small number of simple image features, such as edges or corners, and then formulating mathematical calculation rules. Afterwards, the appropriateness is subjectively determined by inspecting the result on example images and example model parameterizations. If the result is not satisfactory the function is tuned or redesigned from scratch. This heuristic approach relies on the designer’s intuition about a good measure of fitness. Earlier works [12, 14] show that this methodology is erroneous and tedious. This traditional approach is depicted to the top in Figure 2.

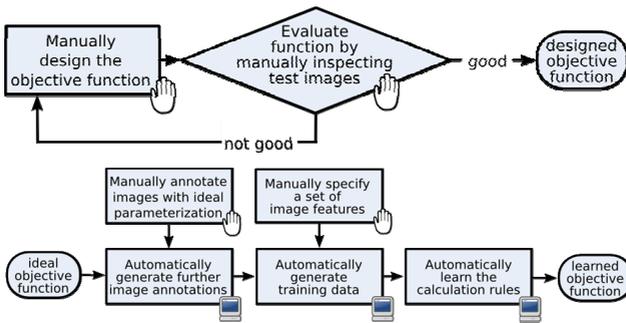


Fig. 2. The traditional procedure for designing objective functions (top), and the proposed method for learning objective functions (bottom)

To avoid these drawbacks, we recently proposed an approach that learns the objective function from annotated example images [14]. It splits up the generation of the objective function into several tasks partly automated. This provides several benefits: firstly, automated steps replace the labor-intensive design of the objective function. Secondly, this approach is less error prone, because giving examples of good fit is much easier than explicitly specifying rules that need to cover all examples. Thirdly, this approach does not rely on expert knowledge and therefore it is generally applicable and not domain-dependent. The bottom line is that this approach yields more robust and accurate objective functions, which greatly facilitate the task of the fitting algorithm.

The fitting algorithm searches for the model that best describes the face visible in the image. Therefore, it aims at finding the model parameters that minimize the objective function. Fitting algorithms have been the subject of intensive research and evaluation, e.g. Simulated Annealing, Genetic Algorithms, Particle Filtering, RANSAC, CONDENSATION, and CCD, see [15] for a recent overview and categorization. We propose to adapt the objective function rather than the fitting algorithm to the specifics

of our application. Therefore, we are able to use any of these standard fitting algorithms, the characteristics of which are well-known, such as termination criteria, runtime, and accuracy.

5 Features Extraction

5.1 Active Shape Models

Different kind of shape models have been introduced by researchers depending upon the application. Some are landmark based models [16,17,18] defining some fixed points annotated on the images and then defining the boundaries around the objects. However some rely on the contour based approach. Different contours define the shape of the object for outlining it along with covering the feature inside an object [19]. Landmark based models however provide the exact location of the features inside the object. Fitting of this shape model on the face is performed by training an objective function (section 4).

The model is parameterized using PCA to form the shape feature vector.

$$x \approx x_m + P_s b_s$$

Where the shape x is parameterized by using mean shape x_m and matrix of eigenvectors P_s to obtain the parameter vector b_s [20].

5.2 Appearance Model

For the various images of the same person different types of variations are required to be modeled. For example, shape deformations including both facial expression changes and pose variations along with the texture variations caused by illuminations. For this reason, different normalizations are required to be performed at this stage. At first, shape variation is required to be controlled in order to record the texture. This can be achieved by defining a reference shape for a particular object. In our case, this reference image is mean shape, obtained by taking the mean of all the shapes of all persons in our database.

Once the texture is extracted it could be parameterized using PCA as,

$$g = g_m + P_g b_g$$

Where the texture g is parameterized by using mean texture g_m and matrix of eigenvectors P_g to obtain the parameter vector b_g [20].

Figure 3 (bottom-left) shows the mean shape of the subject in consideration. Since the points distribution defines a convex hull of points in space so a planar subdivision is defined for the reference shape to map image texture. Delaunay triangulation is used to divide the shape into a set of different facets. Figure 3 shows the delaunay triangulations of the reference shape.

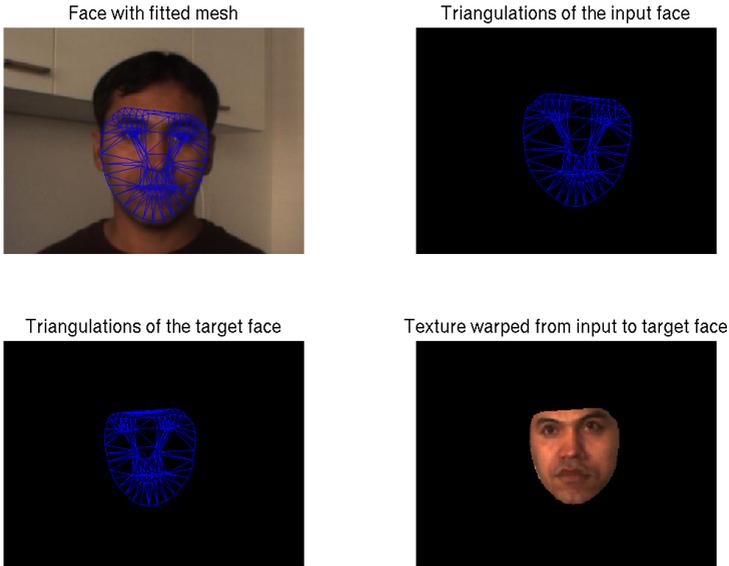


Fig. 3. Shape Subdivisions and Texture Warping

5.3 Temporal Features

Since facial expressions emerge from muscle activity, the motion of particular feature points within the face gives evidence about the facial expression. These features further help the classifier to learn the motion activity. Real-time capability is important, and therefore, a small number of feature points are considered only. The relative location of these points is connected to the structure of the face model. Note that we do not specify these locations manually, because this assumes a good experience of the designer in analyzing facial expressions. In contrast, we automatically generate G feature points that are uniformly distributed. We expect these points to move descriptively and predictably in the case of a particular facial expression. We sum up the motion $g_{x,i}$ and $g_{y,i}$ of each point $1 \leq i \leq G$ during a short time period. We set this period to 2 sec to cover slowly expressed emotions as well. The motion of the feature points is normalized by the affine transformation of the entire face (tx , ty , s , and θ) in order to separate the facial motion from the rigid head motion. In order to determine robust descriptors, PCA determines the H most relevant motion patterns (principal components) visible within the set of training sequences. A linear combination of these motion patterns describes each observation approximately correct. This reduces the number of descriptors ($H \leq 2G$) by enforcing robustness towards outliers as well. As a compromise between accuracy and runtime performance, we set the number of feature points to $G = 140$ and the number of motion patterns b_i to $H = 14$ containing.

The overall feature vector then becomes:

$$u = (b_{s,1}, \dots, b_{s,m}, b_{t,1}, \dots, b_{t,H}, b_{g,1}, \dots, b_{g,n})$$

Where b_s , b_t and b_g are shape, temporal and textural parameters respectively.

6 Feature Classification

With the knowledge of feature the vector u , a classifier infers the correct facial identity. We learn a Binary Decision Tree and Bayesian classifier. However, any other multi-class classifier that is able to derive the class membership from real valued features can be integrated as well, such as a k-Nearest-Neighbour classifier. We take 66% of the image sequences of the CKFE-DB as the training set and the remainder as test set for BDT, the evaluation on which is shown in the next section.

7 Experiments

Experiments have been performed on Cohn-Kanade-Facial-Expression database (CKFE-DB) for human faces. The CKFE-DB contains 488 short image sequences of 97 different persons performing the six universal facial expressions [21]. It provides researchers with a large dataset for experimenting and benchmarking purpose. Each sequence shows a neutral face at the beginning and then develops into the peak expression. Furthermore, a set of action units (AUs) has been manually specified by licensed Facial Expressions Coding System (FACS) [22] experts for each sequence. Note that this database does not contain natural facial expressions, but volunteers were asked to act. Furthermore, the image sequences are taken in a laboratory environment with predefined illumination conditions, solid background and frontal face views. Algorithms that perform well with these image sequences are not immediately appropriate for real-world scenes.

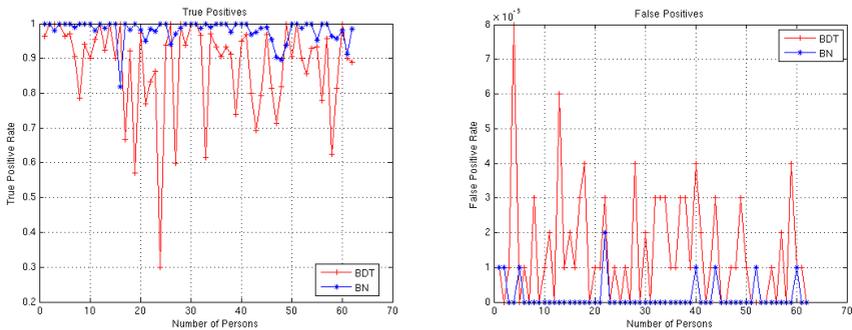


Fig. 4. True Positive and False Positive for 61 persons in experiments

For experimental purposes, image sequences of 62 persons have been used which consists of overall 4060 images. A binary decision tree is trained as classifier in 22.99 sec. We used 1381 images for testing the recognition results and successfully recognized 1259 images. The recognition rate achieved was 91.17% in the presence of facial expressions. Further, same set of the feature vectors are used to build a Bayesian Networks (BN) classifier. The classifier was built in 18.33 sec using 10-fold cross validation. This produces a better accuracy of 98.69% of accurate recognition. Figure 4 shows true positive and true negative for the database for BDT and BN respectively.

8 Conclusions

We introduced an idea to develop a feature vector which consists of three types of facial variations and is robust against the expressional changes in the human faces in real environments. Since the training set consists of the facial expressions information of a person, it can recognize the person even under various expressions. A Bayesian Network (BN) classifier is efficient to train and classify. However the benchmarked database consists of only frontal view of faces. This technique is capable of working in real time environment. This system is a constituent of HRI system. It can keep the person identity information even under the presence of facial expressions which could originate under human machine interaction scenarios. However, in real time environment the system can work by further improving it for light illuminations and using 3D information.

References

- [1] Zhao, W., Chellapa, R., Rosenfeld, A., Philips, P.J.: Face Recognition: A Literature Survey, UMD CFAR Technical Report CAR-TR-948 (2000)
- [2] Barrett, W.A.: A Survey of Face Recognition Algorithms and Testing Results. Proceeding of IEEE (1998)
- [3] Sirovich, L., Kirby, M.: Low-dimensional procedure for the characterization of human faces. *J. Opt. Soc. Am. A* 4(3), 519–524 (1987)
- [4] Turk, M.A., Pentland, A.P.: Face Recognition using Eigenfaces. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 586–591 (1991)
- [5] Turk, M.A., Pentland, A.P.: Eigenfaces for Recognition. *Journal of Cognitive Neuroscience* 3(1), 71–86 (1991)
- [6] Edwards, G.J., Taylor, C.J., Cootes, T.F.: Interpreting Face Images using Active Appearance Models. In: Proceedings of International Conference on Automatic Face and Gesture Recognition, pp. 300–305 (1998)
- [7] Edwards, G.J., Cootes, T.F., Taylor, C.J.: Face Recognition using Active Appearance Models. In: Burkhardt, H., Neumann, B. (eds.) ECCV 1998. LNCS, vol. 1407, pp. 581–695. Springer, Heidelberg (1998)
- [8] Edwards, G.J., Lanitis, A., Taylor, C.J., Cootes, T.: Statistical Models of Face Images: Improving Specificity. In: British Machine Vision Conference 1996, Edinburgh, UK (1996)

- [9] Belhumeur, P.N., Hespanha, J.P., Kriegman, D.J.: Eigenfaces vs Fisherfaces: Recognition using Class Specific Linear Projection. *IEEE Transaction on Pattern Analysis and Machine Intelligence* 19(7) (July 1997)
- [10] Wimmer, M., Riaz, Z., Mayer, C., Radig, B.: Recognizing Facial Expressions Using Model-Based Image Interpretation. In: *Advances in Human-Computer Interaction*, vol. 1, pp. 587–600 (I-Tech Book)
- [11] Lantis, et al.: A unified Approach to Coding and Interpreting Face Images. In: *5th International Conference on Computer Vision*, Cambridge USA (1995)
- [12] Wimmer, M., Stulp, F., Tschechne, S., Radig, B.: Learning Robust Objective Functions for Model Fitting in Image Understanding Applications. In: *Proceedings of the 17th British Machine Vision Conference, BMVA*, Edinburgh, UK, pp. 1159–1168 (2006)
- [13] Cootes, T.F., Taylor, C.J.: Active Shape models – smart snakes. In: *Proceedings of the 3rd British Machine Vision Conference*, pp. 266–275. Springer, Heidelberg (1992)
- [14] Wimmer, M., Stulp, F., Pietzsch, S., Radig, B.: Learning local objective functions for robust face model fitting. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* 30(8) (2008)
- [15] Hanek, R.: Fitting Parametric Curve Models to Images Using Local Self-adapting Separation Criteria. PhD thesis, Department of Informatics, Technische Universität München (2004)
- [16] Cootes, T.F., Edwards, G.J., Taylor, C.J.: Active Appearance Models. In: Burkhardt, H., Neumann, B. (eds.) *ECCV 1998. LNCS*, vol. 1407, pp. 484–498. Springer, Heidelberg (1998)
- [17] Stegmann, M.B.: Active Appearance Models: Theory Extensions and Cases, Master Thesis, Technical University of Denmark (2000)
- [18] Ahlberg, J.: An Experiment on 3D Face Model Adaptation using the Active Appearance Algorithm, Image Coding Group, Dept. of Electric Engineering, Linköping University
- [19] Gupta, H.R., Chowdhury, A.K., Chellappa, R.: Contour-based 3D Face Modeling from a Monocular Video. In: *British Machine Vision Conference 2004* (2004)
- [20] Li, S.Z., Jain, A.K.: *Handbook of Face recognition*. Springer, Heidelberg (2005)
- [21] Kanade, T., Cohn, J.F., Tian, Y.: Comprehensive database for facial expression analysis. In: *Proceedings of Fourth IEEE International Conference on Automatic Face and Gesture Recognition (FG 2000)*, Grenoble, France, pp. 46–53 (2000)
- [22] Ekman, P., Friesen, W.V.: *The Facial Action Coding System: A Technique for the Measurement of Facial Movement*. Consulting Psychologists Press, San Francisco (1978)
- [23] Ross, A.: *Introduction to Biometrics*. Lecture in 5th International Summer School of Biometrics, Italy (June 2008)