

Face Recognition Using Ensembles of Networks

S. Gutta⁺, J. Huang⁺, B. Takacs^{*}, and H. Wechsler⁺

⁺ Department of Computer Science

^{*} Computational Sciences and Informatics

George Mason University

Fairfax, VA 22030

Abstract

We describe a novel approach for fully automated face recognition and show its feasibility on a large data base of facial images (FERET). Our approach, based on a hybrid architecture consisting of an ensemble of connectionist networks - radial basis functions (RBF) - and inductive decision trees (DT), combines the merits of 'discrete and abstractive' features with those of 'holistic' template matching'. Training for face detection takes place over both positive and negative examples. The benefits of our architecture include (i) robust detection of facial landmarks using decision trees, and (ii) robust face recognition using consensus methods over ensembles of RBF networks. Experiments carried out using k - fold cross validation on a large data base consisting of 748 images corresponding to 374 subjects, among them 11 duplicates, yield on the average 87 % correct match, and (ROC curves where) 99 % correct verification is achieved for a 2 % reject rate.

1. Introduction

Face recognition is a difficult task because of the inherent variability of the image formation process in terms of image quality and photometry, geometry, occlusion, change, and disguise. Two recent surveys on face recognition discuss these challenges in some detail [1][2]. All face processing systems available today can only perform on restricted data bases of images in terms of size, age, gender, and/or race, and they further assume well controlled environments. There are additional degrees of variability ranging from those assuming that the position/cropping of the face and its environment (distance and illumination) are totally controlled, to those involving little or no control over the background and viewpoint, and eventually to those allowing for major changes in facial appearance due to factors such as aging and disguise (hat and/or glasses).

For most forensic scenarios there are two possible recognition tasks to be considered:

- MATCH: An image of an unknown individual is collected and the identity has to be found by searching a large set of images.

- SURVEILLANCE: The task is that of verification - whether a given probe does belong to a relatively small gallery, possibly consisting of just one image.

There are two major approaches for automated identification of human faces. The first approach, the *abstractive* one, extracts (and measures) discrete local features 'indexes' for retrieving and identifying faces, while standard statistical pattern recognition techniques are then employed for matching faces using these measurements. Examples of the feature extraction step include the discrete Karhunen-Loeve (KL) and the measurement of anthropomorphic characteristics [3]. The other approach, the *holistic* one, conceptually related to template matching, attempts to identify faces using global representations. Characteristic of this approach are connectionist methods such as backpropagation using holons [4], principal component analysis (PCA), and singular value decomposition (SVD) using eigenfaces [5].

The approach described in this paper is hybrid in its nature and it is based on a psychologically plausible distinction between two types of cognitive operations: automatic, reflexive or low level (e.g., perception) vs. controlled, deliberative or high level (e.g., reasoning). The concept of reductionism is a common practice in the development of intelligent systems - to design solutions to complex problems through a stepwise decomposition of the task into successive modules. In the context of face recognition, hybrid architectures, consisting of connectionist networks and symbolic methods, would thus combine the merits of 'discrete' methods using numerical and symbolic values, with those of 'holistic' template matching, respectively. The hybrid approach for classification involves specific levels of knowledge where the hierarchy is defined in terms of concept granularity. As one moves upward in the hierarchical structure, we witness a corresponding degree of data compression so more powerful ('reasoning') methods can be employed but on reduced amounts of data.

We describe in this paper automated solutions for both the MATCH and SURVEILLANCE tasks, and present experimental results for large facial (image database) test beds as those made available under the FERET research program. The FERET facial data base is described in Sect. 2, learning from examples in Sect. 3, ensembles of networks and consensus methods are

discussed in Sect. 4, while the methodology used for face recognition is presented in Sect. 5. The Automated Face Recognition (AFR) system, its implementation, and experimental results are described in Sect. 6. Conclusions and future enhancements are discussed in Sect. 7.

2. The FERET Facial Database

For the most part, the performance of face recognition systems reported in the literature has been measured on small databases, with each research site carrying out its experiments on their own database thus making meaningful comparisons and drawing conclusions impossible [6]. The majority of those databases were collected under very controlled situations and algorithms were developed to process mostly frontal or full profile images. To overcome such shortcomings, we have been developing over the last several years the FERET facial database so a standard tested for face recognition applications can become available [7,8]. The FERET data base consists now of 1, 109 sets comprising 8, 525 images. Since large amounts of images were acquired during different photo sessions, the lighting conditions and the size of the facial images can vary. The diversity of the FERET data base is across gender, race, and age. Fig.1 is indicative of the range of facial images the FERET DB consists now of and it will be referred to when we describe different stages of automatic face recognition.



Figure 1. Examples of Facial Images

1,109 facial image sets - including 190 duplicate sets taken at different times and possibly wearing glasses - consisting of several poses and totaling 8,525 images have been collected so far. Acquisition of duplicate sets is very important if one wants to assess how robust is a given face recognition system when tested on images shot at different times, which are likely to be somehow different. Most of the sets consist of the following

poses: two frontal shots ('fa' and 'fb'), 1/4 half (right and left) profiles ('qr' and 'ql'), 3/4 half (right and left) profiles ('hr' and 'hl'), and right and left (90 deg.) profiles ('pr' and 'pl'). This paper is concerned with frontal images only. As we are already exploring how to process unrestricted face orientation, we recently collected several hundred sets that have several additional poses at the midpoints between: 'hr' and 'qr' ('ra'), 'qr' and frontal view ('rb'), frontal view and 'ql'('rc'), 'ql' and 'hl' ('rd'), and between 'hl' and 'pl' ('rd'). The additional poses were taken to assess the capability of modeling the human face using several positions and interpolating/extrapolating amongst them for identification purposes. The facial image sets were acquired without any restrictions imposed on expression and with two frontal images shot at different times during the photo session.

3. Learning from Examples

Inductive Inference is defined as the process of going from specific observational knowledge about some objects and a (possibly null) initial inductive hypothesis to an inductive assertion that "strongly" or "weakly" implies or accounts for the observations. One subdomain of inductive inference is *concept learning from examples*, in which the specific knowledge consists of a set of objects belonging to known classes. The inductive assertion is expressed as a *classification rule* for assigning any object, seen or unseen, to a class. These objects are known only through their *descriptions* in terms of a collection of properties, which might include measurements, yes-no indicators, and qualitative measurements.

Concept learning takes place as a set of examples namely positive and counter positive or negative examples is presented. The basic principle is to learn a description of an concept that primarily covers only the positive examples but none of the negative examples. This notion can also be extended to have more than two classes. One distinct advantage of these systems includes learning of adaptive thresholds so the need for empirical thresholds is eliminated. Characteristic of this approach are the Version Space approach (VS), Algorithm Quasi-optimal (AQ), and Inductive Decision trees (ID3, C4.5) [8].

The methodology described in Sect. 5 and the automated face recognition architecture described in Sect. 6 implement the eye detection stage, necessary for face normalization, using learning from examples. The novelty of using this approach for face recognition in general, and facial landmarks ('eye') detection in particular, stems from us employing both positive and negative examples.

4. Ensembles of Networks and Consensus Methods

An early example of using ensembles of neural networks is due to Hamshire and Waibel [10]. The

Meta - Pi classifier is a connectionist pattern classifier that consists of a number of source-dependent sub networks that are integrated by a combinational Time Delay Neural Network (TDNN) superstructure. The TDNN combines the outputs of the modules, trained independently, in order to provide a global classification. Lincoln and Skrzypek [11] have proposed a clustering multiple backpropagation networks for improved performance and fault tolerance. Following training, a 'cluster' is created by computing the average of the outputs generated by the individual networks. The output of the 'cluster' is used as the desired output during training by feeding it back to the individual networks. The basic notion behind using such a strategy is based according to the authors on the assumption that while it is possible to 'fool' single BP networks all of the time one cannot mislead all of them all of the time.

Battiti and Colla [12] have proposed means to combine the outputs of different neural network classifiers to improve the rejection-accuracy (ROC) rates and to make the combined performance better than that obtained from the individual components. The suggested concept of democracy is analogous with the human way of reaching a pondered decision - query by consensus . Flocchini *et al.* [13], have proposed a complex architecture based on a hierarchy of neural networks with a self-referencing structure. The system is structured as a tree in which nodes correspond to neural networks, each one having different tasks. Each leaf is a recognition module composed by some networks with different characteristics. These networks are coordinated by a supervisor in a self-referencing structure. During the training phase, the Meta-Net supervisor observes the behavior of recognition nets and learns which net is more reliable in what task. During the test phase the Meta-Net decides, given an input image, what weights to assign to each network and how to modify their output in order to obtain the final result.

Connectionist architectures are successful when they can cope with the variability of the image acquisition process. One possible solution to the above problem is to implement the equivalent of query by consensus using ensembles of radial basis functions (ERBF) and to train them on data reflecting the inherent variability of the input. Ensembles are defined in terms of their specific topology (connections and RBF nodes) and the data they are trained on. Specifically, both original data and distortions caused by geometrical changes and blur are used to train so robustness to different distortions is achieved via generalization ('hints') [14]. The specific realization of such ERBF is discussed in Sect. 6.

5. Methodology

The overall architecture is shown in Fig. 2. Face recognition usually starts through the detection of a pattern as a face and boxing it, proceeds then by normalizing the face image to account for geometry and

illumination changes using information about the box surrounding the face and/or the eye location, and finally it identifies the face using appropriate classification algorithms. The tools developed to realize the AFR architecture and implement the face recognition process mentioned above are discussed in detail in Sect. 6.

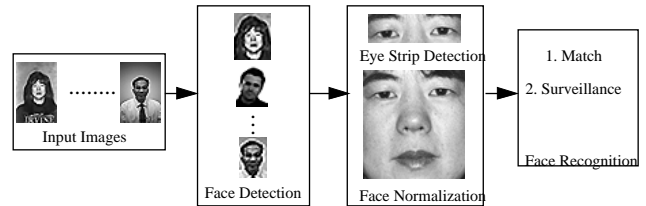


Figure 2. Architecture of the Automated Face Recognition (AFR) System

6. Automated Face Recognition (AFR)

We describe next each processing stage leading to face recognition and we address specific issues related to their computer implementation. The AFR addresses both the MATCH and SURVEILLANCE tasks as defined earlier. We experimented with a database of 748 images drawn randomly from batches 1, 2, 3 and 4. These batches have been acquired at different times, namely, batch 1 on 08/31/93; batch 2 on 12/30/93; batch 3 on 01/28/94 and batch 4 on 04/22/94. The database includes 374 'fa' and 374 'fb' of which 11 pairs of images are duplicates. All images are of size 256 x 384 using 256 gray levels.

6.1 Face Detection

To quickly locate the bounding box of the face (see Fig. 3) before more expensive facial feature (eyepair detection) algorithms are executed, we use a simple algorithm that operates on the edge image. First the edge image at resolution 64 x 96 pixels is created using the MATLAB implementation of the Marr edge detector. Horizontal and vertical projection profiles are computed to find the maximum projection ($P_{h,max}$ and $P_{v,max}$) values that would mark possible regions of interest. After thresholding/filtering the horizontal projection profile, by using only those lines where the projected values are greater than $0.8 \times P_{h,max}$, the upper boundary of the face box is found and it establishes a reference point for further processing. Once the upper boundary is detected, the algorithm searches for the maximal vertical projections within a certain distance from the center line to find the right and left sides of the bounding box. The lower boundary is finally computed using a fixed ratio of the upper boundary and the two edges on each side so the ratio height (h_0) / width (w_0) = 1.5. The accuracy of this stage, determined using visual inspection, is 94.5 %, and it corresponds to 42 incorrect box images (See Table 1).

An incorrect box means that not all the facial features are present. The reason for failing to detect the right box stems primarily from low contrast. Note that all the images, including those for whom an incorrect box is found, are passed to the next stage, that of face normalization.



Figure 3. Face Detection

6.2 Face Normalization

Face normalization is carried out using the location of the eye strip. As a consequence two substages are implemented here: one would locate the eye strip and the other will use the midpoint of the eye strip for normalization. As it would be discussed later on learning involves both positive and negative instances ('examples').

6.2.1 Eye Strip Detection

Once the face box has been detected we proceed to locate the eye strip so an anchor for future normalization can be derived. The cropped ('boxed') face images are now made available at an increased resolution of 128 x 192, which is twice the resolution of the images used for cropping the face. The reason for affording increased resolution is that we have to seek finer facial details within a restricted area. The strategy for eye strip detection learns to label possible facial sites as eye strips using decision trees (DT) defined in terms of appropriate features. The eye strip detection process, explained in detail below, consists of: data preparation, feature extraction, derivation of optimal decision tree (C4.5), and proper training and testing.

6.2.1.1 Data Preparation

The data set used to train the eye strip detection comes from 40 images and would consist of 40 (correct / positive: '+' exs.) eye-strips of dimension 32 x 48 and 200 (incorrect / negative: '-' exs.) non eye-strip images. The actual eye-strips are usually of dimension 16 x 48 - the reason for using positive examples of larger size is eye strip detection would be enhanced when contextual cues given by the cheeks are taken into consideration. The eye strip is partitioned into a 2 x 3 arrays of windows where the eyes correspond to the left and right windows in the top row and immediately below we have the corresponding cheek windows.

6.2.1.2 Feature Detection

The set of eleven features is based on (a) the Laplacian, (b) ranking the Laplacian values, and (c) local symmetry, and it is described below.

(a) Laplacian

The Laplacian, provides for invariant features not affected by lighting conditions, and it detects image transitions usually associated with edge segments and/or image contrast. The Laplacian operator is defined as

$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

Four features are then computed and correspond to the mean of the Laplacian in each of the subwindows mentioned above except the middle ones.

(b) Ranking of the Laplacian

The leftmost and rightmost windows define now triangular areas whose base is the outer boundary of the strip and whose vertex is the corresponding mid-point across the boundary. The average values of the Laplacian in these triangles are ranked and four additional features in terms of equivalence classes are defined.

(c) Local Symmetry

As symmetry is so important in describing facial landmarks and the face itself, the last three features measure symmetry characteristics. The features are defined across subtriangles derived from those defined in (b) above when split again into two right-angle triangles and measure the difference between the top three left and right subtriangles.

6.2.1.3 Labeling and Post Processing

Inductive learning, as applied to the detection of eye-strips over the face region requires a special interface of numeric-to-symbolic data conversion. This has been achieved by tagging each of the 40 positive examples of eye-strips as 'CORRECT' and each of the 200 negative examples of non-eye-strips as 'INCORRECT'. The inductive learning system that has been used for this purpose is Quinlan's C4.5 [9]. The input to the C4.5 consists of a string of learning events, each event given as a vector of attribute values. C4.5 is an inductive machine learning system, which takes a set of positive examples and a set of negative examples and builds a classifier as a decision tree whose structure consists of

- *leaves*, indicating class identity, or
- *decision nodes* that specify some test to be carried out on a single attribute value, with one branch for each possible outcome of the test.

A decision tree is used to classify an example by starting at the root of the tree and moving through it until a leaf is encountered. At each nonleaf a decision is evaluated, the outcome is determined, and the process moves on accordingly. The decision nodes implement an optimal (entropy) criterion such as the gain ratio criterion [9].

Once training has been completed and the learning system has found a suitable decision tree, specific testing and/or AFR system operation can proceed. This stage is performed by moving a window of size 32 x 48 with an overlap factor of size 6 x 8 pixels over the entire face box image. From each scanned window a set of eleven feature vectors are extracted as explained in

Section 5.2.1.2. For each image, the windows that have been correctly classified are cut out of the face box image. As more than one eye strip is sometimes detected for each face box image, postprocessing such as winner-take-all (WTA) is needed, whereby the set of candidate eye strips is reduced to one strip. This is achieved by using C4.5 and retraining based on the relative location of the boxes in the cropped face box region. In cases where more than one eye strip still competes to win simple averaging of the coordinates yields the best location. Examples of eye strip detection are presented in Fig. 4. Note that there is a possibility for the decision tree classifier to fail in detecting eye strips; if this is the case, face recognition stops and the image is not processed any further. This accounts for the images counted as rejected in Table 1 -- 18 images rejected during eye strip labeling and 8 additional images rejected during postprocessing.

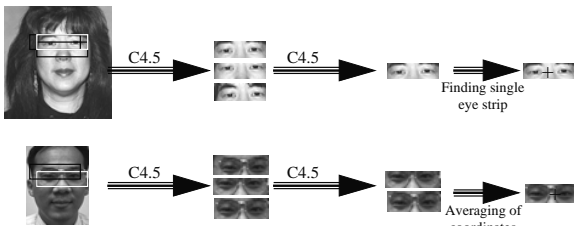


Figure 4. Eye Strip Detection

6.2.1.4 Performance Evaluation

The face detection stage, as mentioned earlier, does not discard any image from further processing. The stage of eye strip detection - labeling and postprocessing, however, would discard images from further processing if it fails to detect the eye strip. Table 1 shows the performance displayed by the AFR system during the face detection, eye strip labeling, and normalization stage. Note that no incorrect decisions are made by the system during the eye strip labeling and postprocessing stage. The net result is that 722 images - 361 pairs - out of 748 images will be passed to the recognition stage, which corresponds to a success rate of 96.5%.

	Stage (section)	Input Size (images)	Processed			Rejected
			Correct	Incorrect	Total	
Face detection	5.1	748	706	42	748	0
Eye strip labeling	5.2.1	748	730	0	730	18
Post processing	5.2.1	730	722	0	722	8
Normalization	5.2.2	722	722	0	722	0

Table 1. Performance Evaluation for Face Detection and Normalization

6.2.2 Face Normalization

Face normalization is performed to ensure that all the face images are of the same size for the classification stage. The eye strip detected in the previous step uniquely defines the anchor point for this normalization stage as the central point of the strip. The coordinates

of the anchor point are mapped at resolution 128 x 192. Average measurements derived during face cropping and the midpoint of the eye strip just detected are used to 'cut' normalized faces for the recognition stage as shown in Fig. 5. The resulting face images are made available to the classification stage at a resolution of 64 x 72.

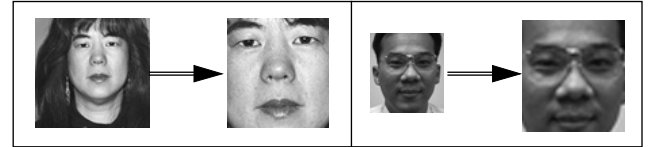


Figure 5. Face Normalization

6.3. Face Recognition

The normalized face images are now ready for recognition. As the amount of information contained in each face image is quite large it seems appropriate that one would first try to cluster the data before its final classification. The Radial Basis Function (RBF) classifier [15] appears to be a natural choice for this task. The reasons behind using RBF are its ability for clustering similar images before classifying them and the potential for developing in the future hierarchical classifiers where faces will be first discriminated in terms of gender, race, and age, before final recognition would take place. Such an approach is akin to that used by people when engaged in face recognition. We start this section by briefly reviewing some of the fundamentals behind RBF and then discuss how the MATCH and SURVEILLANCE tasks are actually performed.

6.3.1 RBF

An RBF classifier has an architecture very similar to that of a traditional three-layer back-propagation network. Connections between the input and middle layers have unit weights and, as a result, do not have to be trained. Nodes in the middle layer, called BF nodes, produce a localized response to the input. That is, each hidden unit can be viewed as a localized receptive field. The hidden layer is trained using k-means clustering. The most common basis function chosen is a Gaussian function, in which the activation level y_i of the hidden unit i is given by:

$$y_i = \Phi_i(\|X - \mu_i\|) = \exp\left[-\sum_{k=1}^D \frac{(x_k - \mu_{ik})^2}{2h\sigma_{ik}^2}\right]$$

where h is a proportionality constant for the variance, x_k is the k th component of the input vector $X=[x_1, x_2, \dots, x_D]$, and μ_{ik} and σ_{ik}^2 are the k th components of the mean and variance vectors, respectively, of basis function node i . The outputs of the hidden unit lie between 0 and 1; the closer the input to the center of the

Gaussian, the larger the response of the node. The activation level Z_j of an output unit is given by:

$$Z_j = \sum_i w_{ij}y_i + w_{0j}$$

where Z_j is the output of the j th output node, y_i is the activation of the i th BF node, w_{ij} is the weight connecting the i th BF node to the j th output node, and w_{0j} is the bias or the threshold of the j th output node. The bias comes from the weights associated with a BF node that has a constant unit output regardless of the input. An unknown vector X is classified as belonging to the class associated with the output node j with the largest output Z_j .

6.3.2 MATCH

The RBF classifier described above has been applied to the MATCH task. Experiments were performed on 361 pairs of images, where each pair has both the frontals 'fa' and 'fb'. Training is performed until either the 'fa' or 'fb' sets are perfectly classified, while testing takes place on the counterpart sets, 'fb' or 'fa', respectively. The first cycle takes a training set consisting of 350 'fa' frontal images, while testing on their corresponding 372 'fb' images - the additional 12 'fb' images correspond to duplicates. The second cycle would change training and testing roles between the 'fa' and 'fb' images. Table 2 below reports the results for the case when correct MATCH corresponds to the best output of RBF or to the first two outputs of RBF.

Cycle	Correct MATCH (1st output)	Correct MATCH (1st output)	Correct MATCH (1st or 2nd output)	Incorrect MATCH (1st or 2nd output)
1	88.89%	11.11%	92.71%	7.29%
2	84.23%	15.77%	89.57%	10.43%
average	86.56%	13.44%	91.14%	8.86%

Table 2. Statistics for the MATCH Task.

1st cycle (fa vs fb)	Correct MATCH	Incorrect MATCH	REJECT Undecided
$\theta = 0.60$	90.34%	9.66%	5.68%
$\theta = 0.65$	93.29%	6.71%	7.34%
$\theta = 0.70$	96.39%	3.61%	12.36%
$\theta = 0.75$	98.79%	1.21%	14.57%
2nd cycle (fb vs fa)	Correct MATCH	Incorrect MATCH	REJECT Undecided
$\theta = 0.60$	93.43%	6.57%	6.21%
$\theta = 0.65$	93.87%	6.13%	6.98%
$\theta = 0.70$	95.79%	4.21%	10.56%
$\theta = 0.75$	97.97%	2.03%	12.38%
average	Correct MATCH	Incorrect MATCH	REJECT Undecided
$\theta = 0.60$	91.88%	8.12%	5.94%
$\theta = 0.65$	93.58%	6.42%	7.16%
$\theta = 0.70$	96.09%	3.91%	11.46%
$\theta = 0.75$	98.38%	1.62%	13.47%

Table 3. Statistics for the MATCH Task with the REJECT (Undecided) Option.

Experiments were carried out to determine the ROC curves when the REJECT / UNDECIDED option becomes available and the results are displayed in Table 3. The REJECT / UNDECIDED option corresponds to that case when the AFR system is undecided about what recognition should be made. Again two cycles are performed and the MATCH is accepted only if its strength is above the threshold θ . Otherwise no decision is made and the image is not classified ('matched'). As one would expect, as the threshold θ is increased the recognition system becomes more conservative in its decisions, more images are rejected, and the percentage of correct MATCHes goes up.

The ROC for the MATCH task with a REJECT (undecided) option is displayed next.

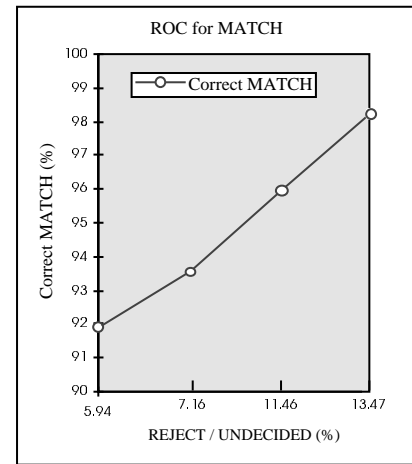


Figure 6. ROC for MATCH with REJECT (Undecided) Option.

Note that the results for duplicate images were consistent for both experiments reported above.

6.3.3 SURVEILLANCE

We also experimented with the same fully automated face recognition system for the SURVEILLANCE task with an REJECT / UNDECIDED option. The experiments were carried out on images drawn randomly from batches 1, 2, 3, and 4 using two slightly different models. The major processing stages for the above system again include (i) face/head cropping, (ii) eye strip pair detection and normalization, and (iii) classification. The first two stages are similar in their implementation to those used for the MATCH task and they were described earlier.

The classification stage employs ERBFs - Ensemble of Radial Basis Functions using the Radial Basis Function (RBF) neural network classifier as building blocks and the strategy used is similar to that of k - fold cross validation (CV) [16]. In k - fold cross validation, the cases are randomly divided into k mutually exclusive test partitions of approximately equal size. The cases not found in each test partition are independently used for training, and the resulting classifier is tested on the

corresponding test partition. The average error rates over all k partitions is the CV error rate. In our experiments, the first CV cycle on its first iteration takes as training set the first 50 'fa' frontal images, while testing on their corresponding 50 'fb' images and the remaining 311 'fa' and 'fb' images - on the second iteration the training set consists of the corresponding 50 'fb' frontal images, while testing on their corresponding 50 'fa' images and the remaining 311 'fb' and 'fa' images. One can define in a similar fashion another seven CV cycles and the performance reported below is the average over seven cycles.

Training is done, as mentioned above, using ERBFs, and the REJECT / UNDECIDED option is implemented. ERBFs are trained over the original images and their slightly degraded (by Gaussian noise, blur, and/or geometrical transformations - small translation and rotation) corresponding frontal view images. The first model, ERBF1, consists of nine RBF networks, each slightly different from the other in terms of the number of clusters and overlap factors. The first three RBF networks are trained over the original images, the next three networks are trained over 'blurred' images, while the last three networks are trained over geometrically distorted images. The final decision as to whether to accept an image or to reject it is based on the following rule - *If the norm of the average of all (nine) outputs is greater than θ accept ('recognize') ; for $(\theta - \Delta\theta, \theta)$ undecided ; else do not accept .* Average cross validation results corresponding to $\theta = 0.67$ are shown below, in Table 4, for different intervals $\Delta\theta$. Specifically, if the average of the ERBF output is within the interval $(\theta - \Delta\theta, \theta)$ no commitment is made as to the identity of the input and the image is labeled as REJECT / UNDECIDED.

av (CV)	Correct SURVEILLANCE	Incorrect SURVEILLANCE	REJECT (Undecided)
$\Delta\theta = 0.02$	97.24%	2.76%	2.75%
$\Delta\theta = 0.06$	97.96%	2.04%	5.14%
$\Delta\theta = 0.10$	98.89%	1.11%	7.61%
$\Delta\theta = 0.14$	99.89%	0.11%	10.22%

Table 4. Statistics for the SURVEILLANCE Task Using the ERBF1 Model.

The ROC curve (correct SURVEILLANCE vs REJECT / UNDECIDED) is displayed in Fig.7.

For the second model, ERBF2, the number of RBF networks was reduced from nine to three by passing the original images and slightly degraded (by Gaussian noise, blur, and/or geometrical transformations - small translation and rotation) corresponding frontal view images to the same net. Each one of the three nets is different from the other in terms of clusters and overlap factors. Though the number of nets is reduced from nine to three, the number of output vectors is still nine. This is obtained by passing the original images and its corresponding slightly degraded images to each one of the three nets, thus generating a total of nine output

vectors. The final decision as to whether to accept an image or to reject it is based on the same rule as used for the first model. Average cross validation results are shown below, in Table 5, for different intervals $\Delta\theta$.

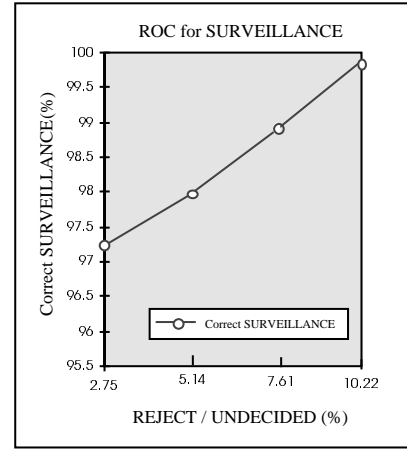


Figure 7. ROC for SURVEILLANCE Using the ERBF1 Model.

av (CV)	Correct SURVEILLANCE	Incorrect SURVEILLANCE	REJECT (Undecided)
$\Delta\theta = 0.02$	98.68%	1.32%	0.95%
$\Delta\theta = 0.06$	99.01%	0.99%	2.00%
$\Delta\theta = 0.10$	99.43%	0.57%	3.51%
$\Delta\theta = 0.14$	99.64%	0.36%	4.76%

Table 5. Statistics for the SURVEILLANCE Task Using the ERBF2 Model.

The ROC curve (correct SURVEILLANCE vs. UNDECIDED) is displayed in Fig. 8.

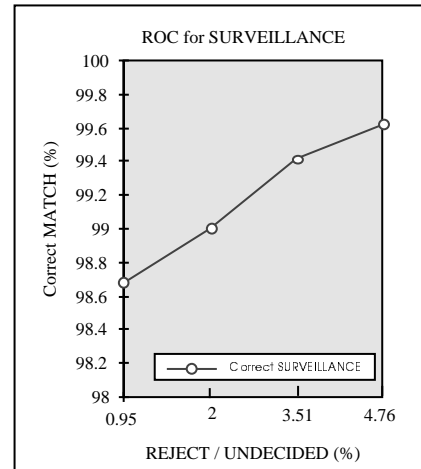


Figure 8. ROC for SURVEILLANCE Using the ERBF2 Model

The rationale for the ERBF models comes from the need for robust performance. Specifically, one should expect that the range for test images is (slightly) different from those encountered during training and

that using more but slightly different nets ('referees') adds to the strength of the decision. The second model, ERBF2, performs slightly better than ERBF1 possibly because of increased flexibility when more than one image per class is made available to each net. Future use of hybrid nets, as used for the eye strip detection, would make the need for empirical thresholds as those used by both ERBFs unnecessary, get rid of the REJECT / UNDECIDED option, and make the whole approach much more robust. The concept of combining the outputs of several neural network classifiers to reach a combined decision with a higher performance, in terms of lower rejection rates and/or better accuracy rates, has been also suggested recently by Battiti and Colla [12] under the label of 'democracy' in neural nets.

7. Conclusions

We described a novel approach for fully automated face recognition and showed its feasibility using a large data base of facial images (FERET). Our approach, based on a hybrid architecture consisting of an ensemble of connectionist networks - radial basis functions (RBF) - and inductive decision trees (DT), combines the merits of 'discrete and abstractive' features with those of 'holistic' template matching'. Training for face detection takes place over both positive and negative examples. The benefits of our architecture include (i) robust detection of facial landmarks using decision trees, and (ii) robust face recognition using consensus methods over ensembles of RBF networks. Experiments carried out using k - fold cross validation on a large data base consisting of 748 images corresponding to 374 subjects, among them 11 duplicates, yield on the average 87 % correct match, and (ROC curves where) 99 % correct verification is achieved for a 2 % reject rate.

We are presently expanding on our work by addressing the issue of how to recognize (for MATCH and SURVEILLANCE) both frontal and slightly off frontal (+/- 15 deg) images, and also images captured at different scales.

Acknowledgements

This work was partly supported by the US Army Research Lab under contract DAAL01-93-K-0099.

References

- [1] Niblack W., "The QBIC project: Querying images by content using color, texture, and shape", *RJ 9203, IBM San Jose Research Division*, 1993.
- [2] Jain A., "Final Report to NSF of the Workshop on Visual Information Management", 1992.
- [3] Craw I.D, Tock B. and A. Bennet, "Finding Face Features", *European Conference on Computer Vision (ECCV)*, Genova, Italy, June 1992.
- [4] Cottrell G.W. and Metcalfe J., "EMPATH: face, gender and emotion recognition using holons", in *Advances in Neural Information Processing Systems (NIPS)*, Vol. 3,

- 564-571, R. P. Lippmann, R.P., Moody J. and Touretzky, D.S. (Eds.), Morgan Kaufmann, 1991.
- [5] Turk M., and Pentland A., "Eigenfaces for Recognition", *J. of Cognitive Neuroscience*, Vol.3, No.1, pp. 71-86, 1991.
- [6] Robertson, G., and I. Craw, "Testing Face Recognition Systems", *Image and Vision Computing*, vol. 19, no. 9, pp. 609-614, 1994.
- [7] DePersia A.T. and P.J. Phillips, "The FERET Program: Overview and Accomplishments", 1995.
- [8] Gutta, S., J. Huang, D. Singh, I. Shah, B. Takacs, and H. Wechsler, "Benchmark Studies on face Recognition", *Proceedings of International Workshop on Automatic Face - and Gesture Recognition (IWAAGR)*, Zurich, Switzerland, June 1995.
- [9] Quinlan, J.R., "The Effect of Noise on Concept Learning", *Machine Learning: an Artificial Intelligence Approach*, R.S. Michalski, J.G. Carbonell and T.M. Mitchell (Eds.), Morgan Kaufmann, pp. 149-166, 1986.
- [10] Hampshire, J. B. and Waibel, A., "The Meta-Pi Network: Building Distributed Knowledge Representations for Robust Multisource Pattern Recognition", *IEEE Trans. on Pattern Analysis and Machine Intelligence* 14, No. 7, 751-769, 1992.
- [11] Lincoln, W. P. and Skrzypek, J., "Synergy of Clustering Multiple Back Propagation Networks", in *Advances in neural Information Processing Systems (NIPS)*, Vol. 2, 650-657, Touretzky, D.S., (Ed.), Morgan Kaufmann, 1990.
- [12] Battiti, R. and Colla, A. M., "Democracy in Neural Nets: Voting Schemes for Classification", *Neural Networks* 7, No. 4, 691-707, 1994.
- [13] Flocchini, P., Gardin, F., Mauri, G., Pensini, M.P. and Stofella, P., "Combining Image Processing Operators and Neural Networks in a Face Recognition System", *International Journal of Pattern Recognition and Artificial Intelligence (IJPRAI)* 6, 446-467, 1992.
- [14] Abu-Mostafa Y. S., "Machines that learn from hints", *Scientific American*, 64-69, 1995.
- [15] Lippmann, R. P., and Ng, K., "A Comparative Study of the Practical Characteristic of Neural Networks and Pattern Classifiers", Technical Report 894, Lincoln Labs., MIT, 1991.
- [16] Weiss, S. M. and C. A. Kulikowski, "Computer Systems That Learn", Morgan Kaufmann, 1991.