

TIED FACTOR ANALYSIS FOR FACE RECOGNITION ACROSS LARGE POSE DIFFERENCES

SIMON J.D. PRINCE, JAMES H. ELDER,
JONATHAN WARRELL, FATIMA M.
FELISBERTI

*IEEE TRANSACTIONS ON PATTERN ANALYSIS AND
MACHINE INTELLIGENCE, JUNE 2008*

Presented by Lan Du
July 25th, 2008

OUTLINE

- Motivation
- Overview of Some Existing Methods for Face Recognition across Pose and the Proposed Method
- Detailed Model and Its Application
 - Observation and Identity Spaces
 - Tied Factor Analysis
 - Learning System Parameters
 - Learning Results
 - Recognition
- Experiments
- Conclusions



MOTIVATION

One of the greatest remaining research challenges in face recognition is to recognize faces across **different poses, expressions, and illuminations**. Current face recognition systems require the implicit cooperation of the user.

- Face recognition from security footage.
- Face recognition in archive footage.
- Face recognition for HCI and ambient intelligence.

In this paper, the authors try to examine the worst case, in which there is **only a single instance of each individual** in a large database, and the *probe* image is taken from **a very different pose** than the matching *gallery* images.



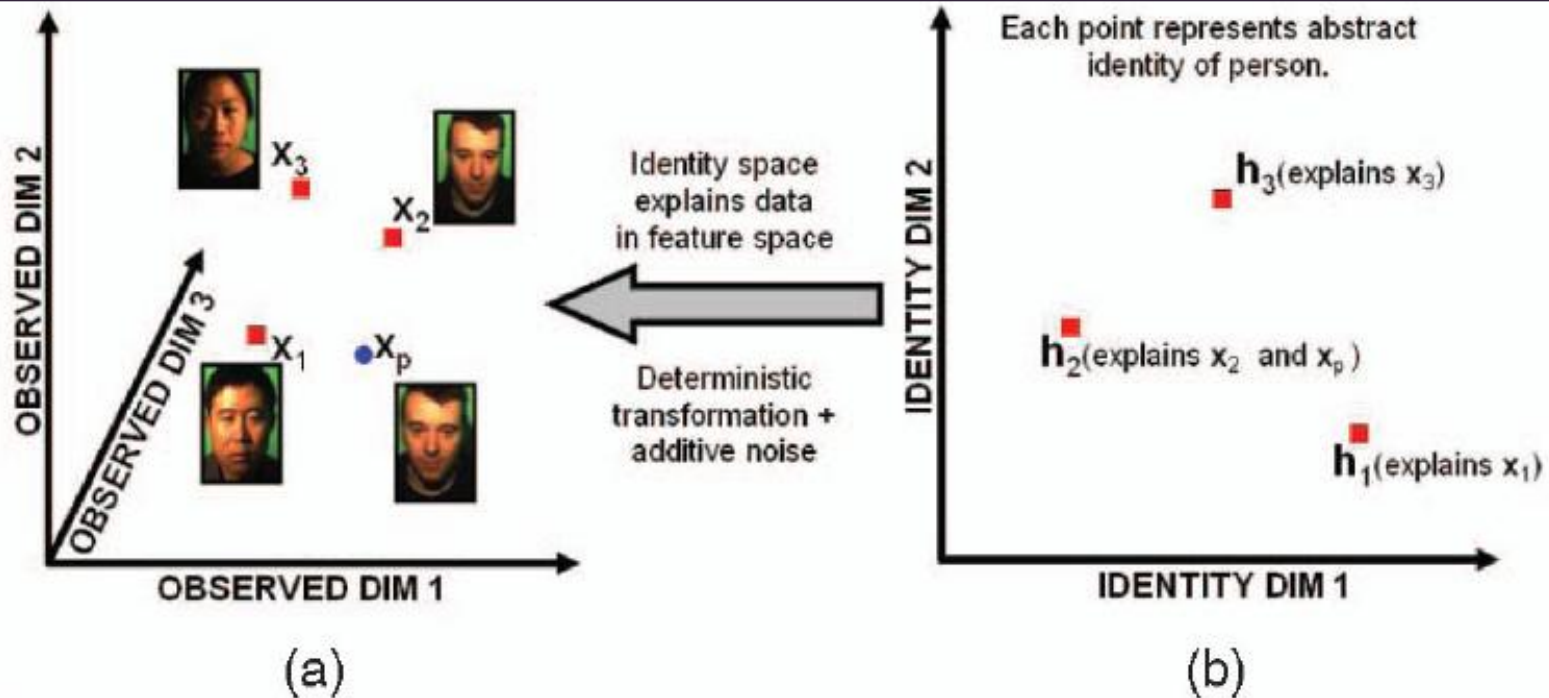
ALGORITHMS FOR FACE RECOGNITION ACROSS POSE

- Record each subject at each possible angle, then use a statistical model for each or create a 3D model of the head. -- require the cooperation of the user
- **3D Geometric Approaches:** take a single probe image at one pose and create a full 3D head. -- complex to implement and are computationally expensive
- **Statistical Approaches:** the relationship between frontal and nonfrontal images is treated as a statistical learning problem. -- simpler and computationally cheaper but produce relatively poor results
 - *Global statistical models*
 - *Local statistical models:* build several models relating different parts of the face



OVERVIEW OF THE PROPOSED METHOD

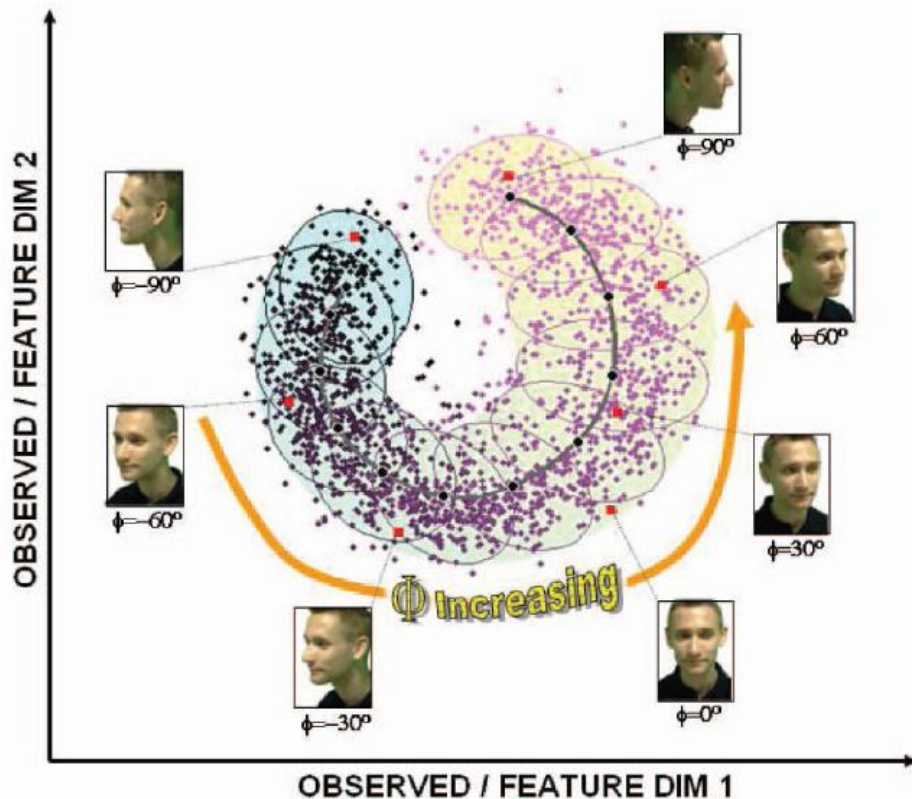
The algorithm is based on a **generative model** that describes how an **underlying pose-invariant representation** created the (pose-varying) observed data.



The latent identity variable approach. (a) Three gallery faces (square symbols) and a probe face (circular symbol) represented in multivariate **observation space**. Each position in this space represents a different image. (b) The "**identity space**", in which each position depicts a different individual. Each image in (a) is modeled as having been generated from a particular point in the identity space in (b).

Observation and Identity Spaces

- **Observed Data:** the raw gray values of the image or some simple deterministic transformation of these values, which does not attempt to compensate for pose variations. -- **Observation Space**
- **Latent Identity Variable:** a multidimensional variable that represents the identity of the individual, regardless of the pose. -- **Identity Space**



The effect of pose variation in the observation space. First, the mean position in the manifold changes systematically with the pose of the face. Second, for a given individual at a given pose, the position of the observation vector, relative to this mean, also varies.



OBSERVATION AND IDENTITY SPACES (CONTD.)

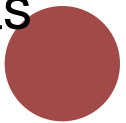
○ Generation Process:

- Choose the point in the identity space that corresponds to an individual.
- Choose a pose.
- Transform this identity variable to the observation space by using a deterministic function, which depends on the pose.
- Add noise to the resulting observation vector.

○ Latent Identity Variable -- describing the shape and structure of the face

Deterministic Function -- representing the perspective projection process, which is parameterized by pose.

Noise Term -- representing the measurement noise in the camera, plus all unmodeled aspects of the situation such as expression and lighting variation.



Tied Factor Analysis

Standard Factor Analysis

$$\mathbf{x}_{ij} = \mathbf{F}_i \mathbf{h}_{ij} + \mathbf{m}_i + \epsilon_{ij}$$

$$\mathbf{h}_{ij} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon_{ij} \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

$$\mathbf{x}_{ij} | \mathbf{h}_{ij} \sim \mathcal{N}(\mathbf{F}_i \mathbf{h}_{ij} + \mathbf{m}_i, \Sigma)$$

$$\begin{aligned} p(\mathbf{x}_{ij} | \mathbf{F}_i, \mathbf{m}_i, \Sigma) \\ &= \int p(\mathbf{h}_{ij}) p(\mathbf{x}_{ij} | \mathbf{h}_{ij}, \mathbf{F}_i, \mathbf{m}_i, \Sigma) d\mathbf{h}_{ij} \\ &= \mathcal{N}(\mathbf{x}_{ij}; \mathbf{m}_i, \mathbf{F}_i \mathbf{F}_i^T + \Sigma) \end{aligned}$$

\mathbf{x}_{ij} : the j th observed data vector of the i th individual

\mathbf{h}_{ij} : factor corresponding to \mathbf{x}_{ij}

\mathbf{F}_i : factor loading for the i th individual

\mathbf{m}_i : offset of the i th individual

ϵ_{ij} : noise term corresponding to \mathbf{x}_{ij}

Tied Factor Analysis

$$\mathbf{x}_{ijk} = \mathbf{F}_k \mathbf{h}_i + \mathbf{m}_k + \epsilon_{ijk}$$

$$\mathbf{h}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon_{ijk} \sim \mathcal{N}(\mathbf{0}, \Sigma_k)$$

$$\mathbf{x}_{ijk} | \mathbf{h}_i \sim \mathcal{N}(\mathbf{F}_k \mathbf{h}_i + \mathbf{m}_k, \Sigma_k)$$

$$\begin{aligned} p(\mathbf{x}_{ijk} | \mathbf{F}_k, \mathbf{m}_k, \Sigma_k) \\ &= \int p(\mathbf{h}_i) p(\mathbf{x}_{ijk} | \mathbf{h}_i, \mathbf{F}_k, \mathbf{m}_k, \Sigma_k) d\mathbf{h}_i \\ &= \mathcal{N}(\mathbf{x}_{ijk}; \mathbf{m}_k, \mathbf{F}_k \mathbf{F}_k^T + \Sigma_k) \end{aligned}$$

\mathbf{x}_{ijk} : the j th observed data vector of the i th individual in the k th pose

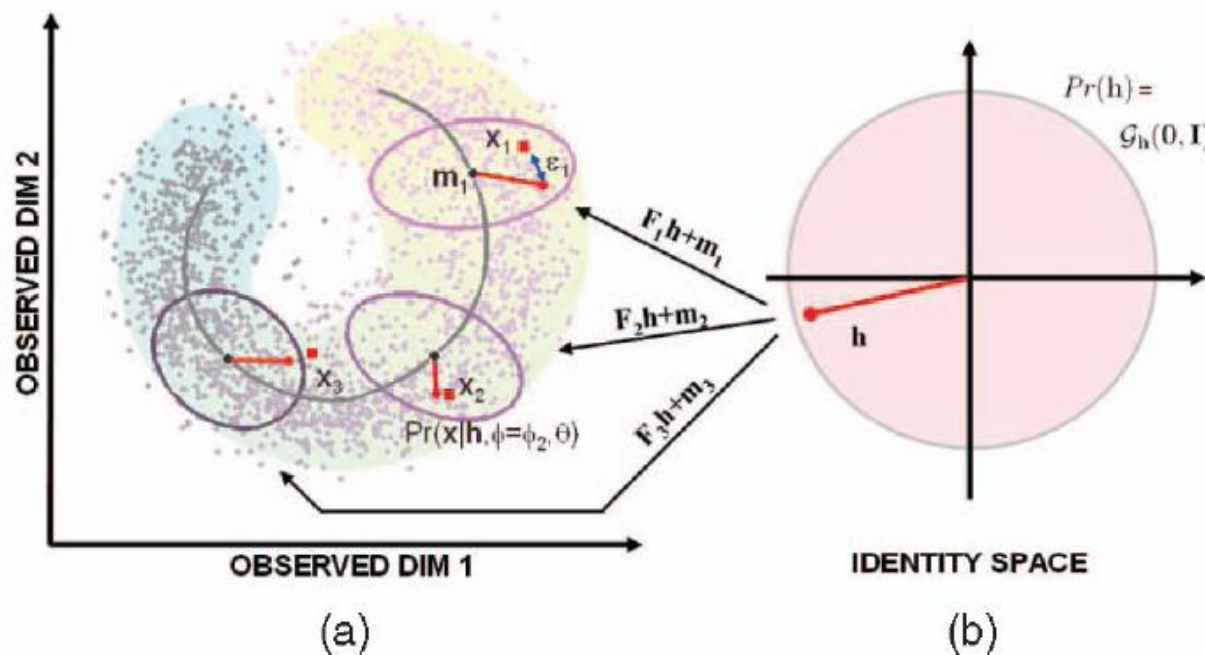
\mathbf{h}_i : factor for the i th individual, which is same at each pose (**tied**)

\mathbf{F}_k : factor loading of the k th pose

\mathbf{m}_k : offset of the k th pose

ϵ_{ijk} : noise term corresponding to \mathbf{x}_{ijk}

TIED FACTOR ANALYSIS (CONTD.)



Tied factor analysis model. (a) Observed measurement space. (b) “Identity” space. The three square symbols in (a) represent observed data for one person viewed at three poses. The circle symbol in (b) represents the latent identity variable for this person. Data in the observation space are explained by transforming latent identity variable by a pose-dependent transform and by adding noise.

LEARNING SYSTEM PARAMETERS

E-Step:

$$\mathbf{E}[\mathbf{h}_i | \mathbf{x}_{i..}] = \left(\mathbf{I} + \sum_{j=1}^J \sum_{k=1}^K \mathbf{F}_k^T \boldsymbol{\Sigma}_k^{-1} \mathbf{F}_k \right)^{-1} \sum_{j=1}^J \sum_{k=1}^K \mathbf{F}_k^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_{ijk} - \mathbf{m}_k)$$

$$\mathbf{E}[\mathbf{h}_i \mathbf{h}_i^T | \mathbf{x}_{i..}] = \left(\mathbf{I} + \sum_{j=1}^J \sum_{k=1}^K \mathbf{F}_k^T \boldsymbol{\Sigma}_k^{-1} \mathbf{F}_k \right)^{-1} + \mathbf{E}[\mathbf{h}_i | \mathbf{x}_{i..}] \mathbf{E}[\mathbf{h}_i | \mathbf{x}_{i..}]^T$$

M-Step:

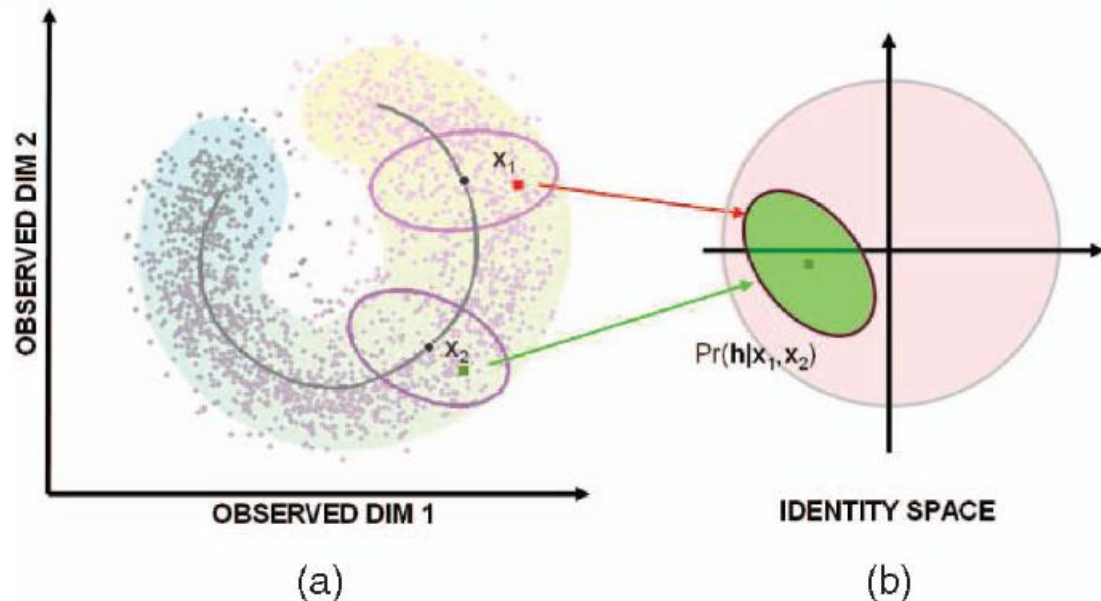
$$\tilde{\mathbf{F}}_k = \left(\sum_{i=1}^I \sum_{j=1}^J \mathbf{x}_{ijk} \mathbf{E}[\tilde{\mathbf{h}}_i | \mathbf{x}_{i..}]^T \right) \left(\sum_{i=1}^I \sum_{j=1}^J \mathbf{E}[\tilde{\mathbf{h}}_i \tilde{\mathbf{h}}_i^T | \mathbf{x}_{i..}] \right)^{-1}$$

$$\tilde{\boldsymbol{\Sigma}}_k = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J \text{diag}[\mathbf{x}_{ijk} \mathbf{x}_{ijk}^T - \tilde{\mathbf{F}}_k \mathbf{E}[\tilde{\mathbf{h}}_i | \mathbf{x}_{i..}] \mathbf{x}_{ijk}^T]$$

$$\tilde{\mathbf{F}}_k = [\mathbf{F}_k \ \mathbf{m}_k], \quad \tilde{\mathbf{h}}_i = [\mathbf{h}_i^T \ 1]^T$$

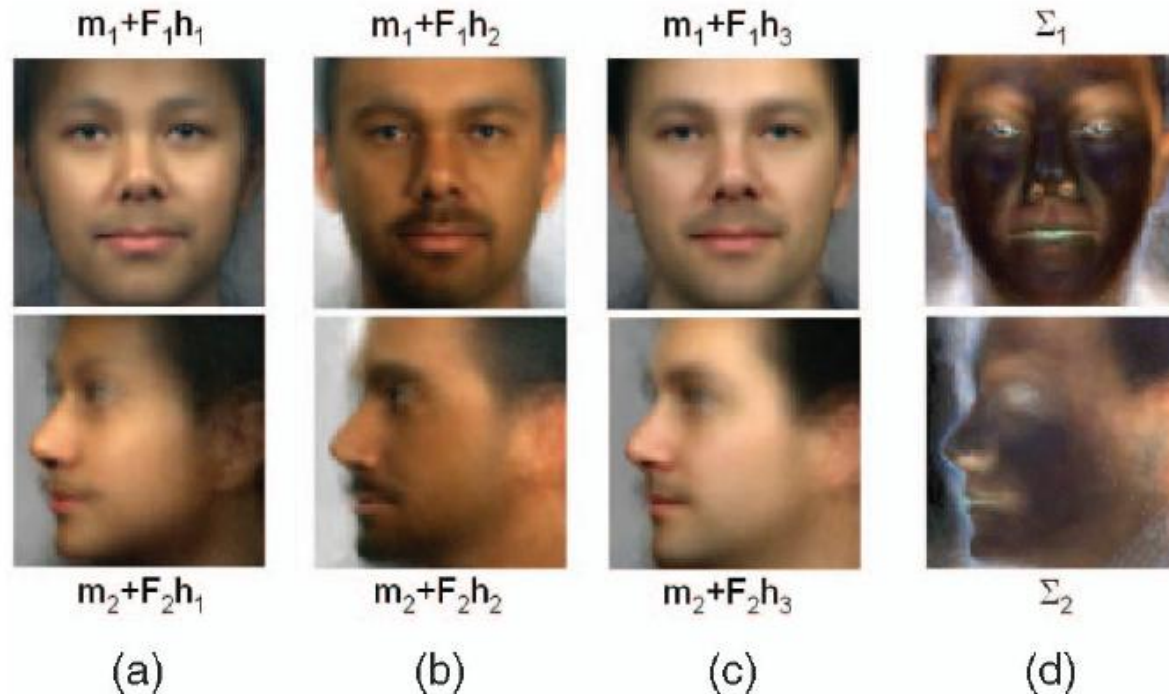
✓(b) The E-Step calculates the posterior probability distribution over the latent identity variables. (a) This is inferred from the observed data for that individual across all poses.

✓The M-Step optimizes the values of the transformation parameters for each pose by using data for that pose across all individuals.



LEARNING RESULTS

- FERET Dataset: 320 individuals at 7 poses -90, -67.5, -22.5, 0, 22.5, 67.5 and 90°; 220 individuals for training and 100 individuals for testing at each pose; identifying 21 keypoints on each face by hand and extracting the corresponding image features.



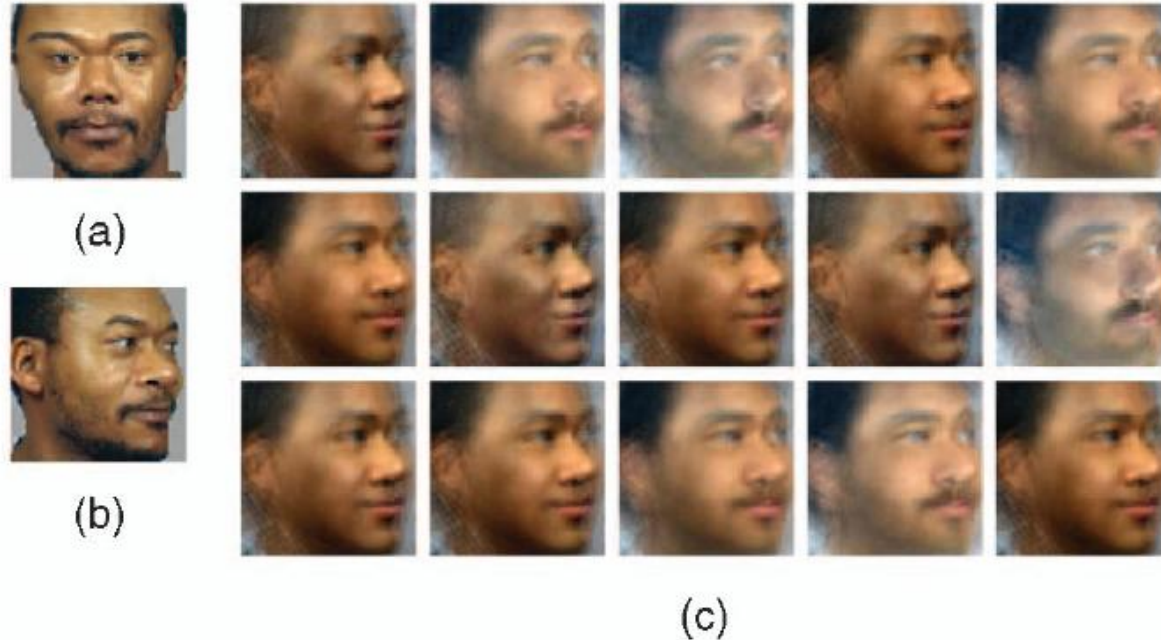
Generated face images with 16 factors. (a), (b), and (c) Three points in the identity space projected back into the observation space through frontal and profile models. (d) Per-pixel noise terms for frontal and profile models. Brighter points represent pixels with more noise.

LEARNING RESULTS (CONT'D.)



Prediction of nonfrontal faces from frontal faces (project the mean of the latent identity variable back to the image space by using a nonfrontal transformation) with 16 factors. (a) Actual images of subject (not in the training database). The frontal image (highlighted in red) is used to predict nonfrontal faces as described in the text. (b) Predicted images for six different poses. (c) (left) One more good example of profile image prediction (left to right: frontal, predicted profile, and actual profile) and (right) one poor example.

LEARNING RESULTS (CONT'D.)



Prediction of nonfrontal faces from frontal faces (project the samples of the latent identity variable back to the image space by using a nonfrontal transformation) with 16 factors. (a) Frontal image of subject. (b) Actual nonfrontal image of subject. (c) Fifteen projected samples.



RECOGNITION

$$\begin{aligned}
 p(\mathbf{x}_{1,\dots,N}, \mathbf{x}_p | \mathcal{M}_n, \boldsymbol{\theta}) &= \int p(\mathbf{x}_{1,\dots,N}, \mathbf{x}_p, \mathbf{h}_{1,\dots,N}, \mathbf{h}_p | \mathbf{h}_p = \mathbf{h}_n, \boldsymbol{\theta}) d\mathbf{h}_{1,\dots,N} \\
 &= \int p(\mathbf{x}_1, \mathbf{h}_1 | \boldsymbol{\theta}_{q_1}) d\mathbf{h}_1 \cdots \int p(\mathbf{x}_n, \mathbf{x}_p, \mathbf{h}_n | \boldsymbol{\theta}_{q_n}, \boldsymbol{\theta}_{q_p}) d\mathbf{h}_n \cdots \int p(\mathbf{x}_N, \mathbf{h}_N | \boldsymbol{\theta}_{q_N}) d\mathbf{h}_N \\
 &= \int p(\mathbf{h}_1) p(\mathbf{x}_1 | \mathbf{h}_1, \boldsymbol{\theta}_{q_1}) d\mathbf{h}_1 \cdots \int p(\mathbf{h}_n) p(\mathbf{x}_n, \mathbf{x}_p | \mathbf{h}_n, \boldsymbol{\theta}_{q_n}, \boldsymbol{\theta}_{q_p}) d\mathbf{h}_n \cdots \int p(\mathbf{h}_N) p(\mathbf{x}_N | \mathbf{h}_N, \boldsymbol{\theta}_{q_N}) d\mathbf{h}_N \\
 &\propto \int p(\mathbf{x}_n, \mathbf{x}_p, \mathbf{h}_n | \boldsymbol{\theta}_{q_n}, \boldsymbol{\theta}_{q_p}) d\mathbf{h}_n / \int p(\mathbf{x}_n, \mathbf{h}_n | \boldsymbol{\theta}_{q_n}) d\mathbf{h}_n
 \end{aligned}$$

$$\int p(\mathbf{x}_{n'}, \mathbf{h}_{n'} | \boldsymbol{\theta}_{q_{n'}}) d\mathbf{h}_{n'} = \mathcal{N}(\mathbf{x}_{n'}; \mathbf{m}_{q_{n'}}, \mathbf{F}_{q_{n'}} \mathbf{F}_{q_{n'}}^T + \boldsymbol{\Sigma}_{q_{n'}}), \quad n' \neq n$$

$$\int p(\mathbf{x}_n, \mathbf{x}_p, \mathbf{h}_n | \boldsymbol{\theta}_{q_n}, \boldsymbol{\theta}_{q_p}) d\mathbf{h}_n = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_n \\ \mathbf{x}_p \end{bmatrix}; \begin{bmatrix} \mathbf{m}_{q_n} \\ \mathbf{m}_{q_p} \end{bmatrix}, \begin{bmatrix} \mathbf{F}_{q_n} \\ \mathbf{F}_{q_p} \end{bmatrix} \begin{bmatrix} \mathbf{F}_{q_n} \\ \mathbf{F}_{q_p} \end{bmatrix}^T + \begin{bmatrix} \boldsymbol{\Sigma}_{q_n} & \\ & \boldsymbol{\Sigma}_{q_p} \end{bmatrix} \right)$$

$$\neq \int p(\mathbf{x}_n, \mathbf{h}_n | \boldsymbol{\theta}_{q_n}) d\mathbf{h}_n \cdot \int p(\mathbf{x}_p, \mathbf{h}_n | \boldsymbol{\theta}_{q_p}) d\mathbf{h}_n$$



$$p(\mathcal{M}_n | \mathbf{x}_{1,\dots,N}, \mathbf{x}_p, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_{1,\dots,N}, \mathbf{x}_p | \mathcal{M}_n, \boldsymbol{\theta}) p(\mathcal{M}_n)}{\sum_{m=1}^N p(\mathbf{x}_{1,\dots,N}, \mathbf{x}_p | \mathcal{M}_m, \boldsymbol{\theta}) p(\mathcal{M}_m)}$$

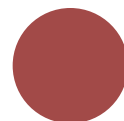
$\mathbf{x}_{1,\dots,N}$: gallery dataset of faces, each of which belongs to a different individual

\mathbf{x}_p : a single probe face

$\boldsymbol{\theta}$: $\{\mathbf{F}_{1,\dots,K}, \mathbf{m}_{1,\dots,K}, \boldsymbol{\Sigma}_{1,\dots,K}\}$

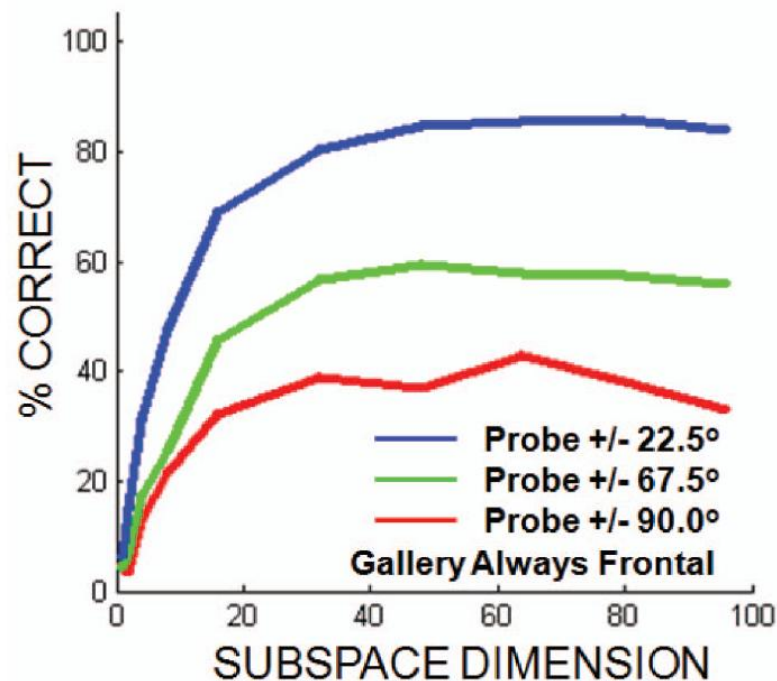
$\boldsymbol{\theta}_{q_n}$: $\{\mathbf{F}_{q_n}, \mathbf{m}_{q_n}, \boldsymbol{\Sigma}_{q_n}\}$

\mathcal{M}_n : the probe matches the n th gallery face

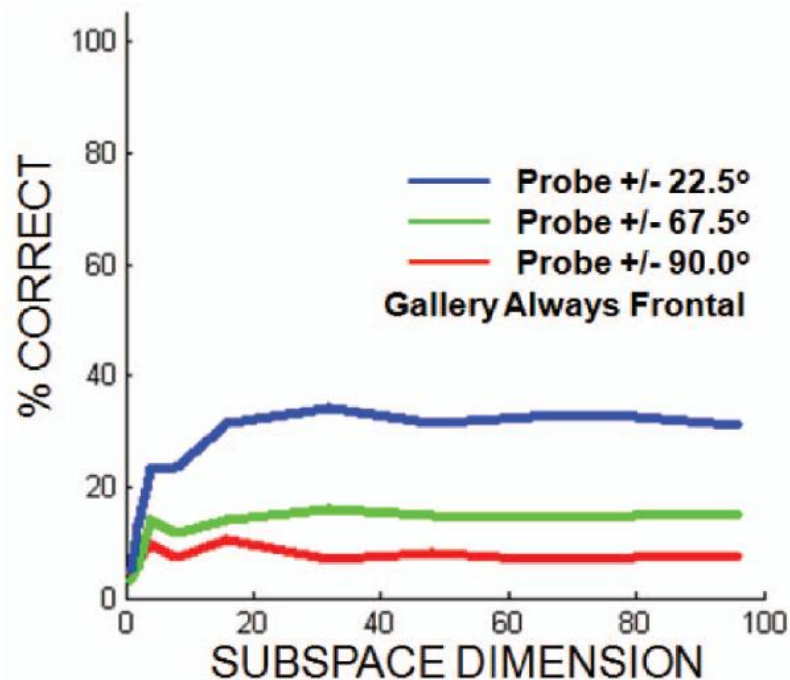


EXPERIMENT 1: FACE IDENTIFICATION USING RAW PIXEL DATA

- 100 frontal testing faces as the gallery faces and a single nonfrontal face as the probe face
- “factor analysis model”: only a single set of generation parameters

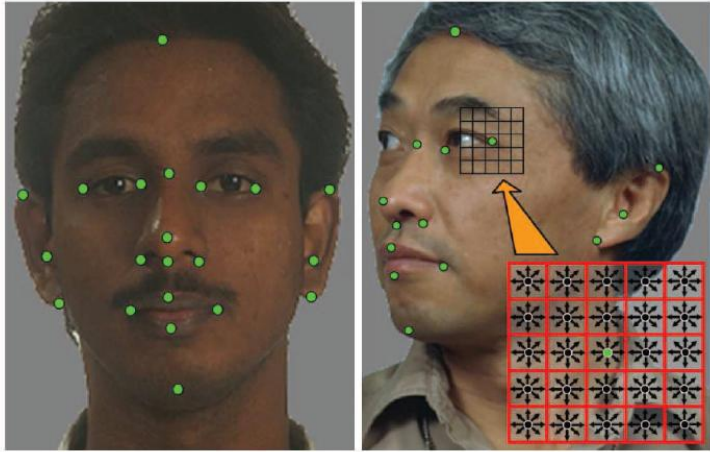


Percentage of first-match correct performance with the tied factor analysis model.



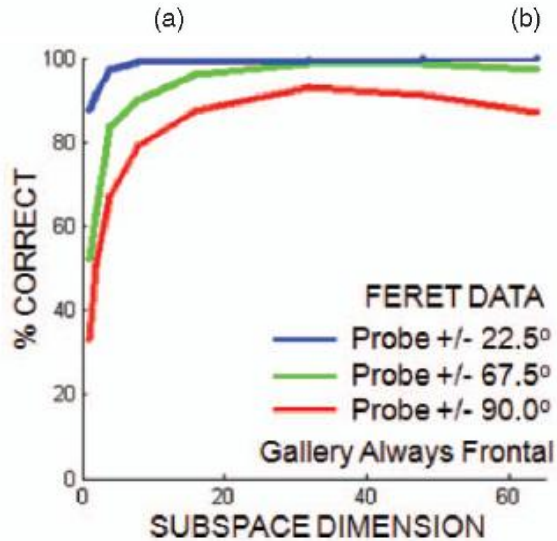
Percentage of first-match correct performance with the “factor analysis model”.

EXPERIMENT 2: FACE IDENTIFICATION WITH LOCAL GABOR DATA

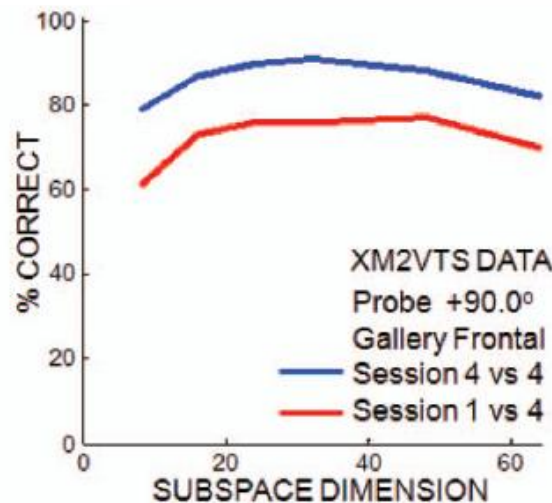


Local measurements. (a) 21 keypoints on each face were identified by hand. (b) features were extracted at 25 spatial positions around each keypoint.

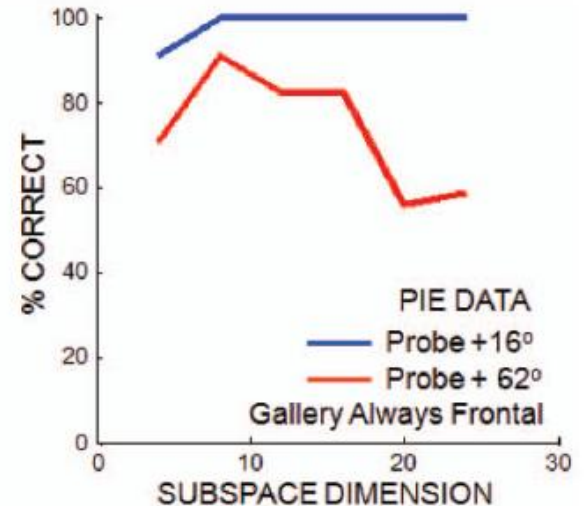
Build 21 **local models** to describe how these local facial features (*nose, eye, etc.*) change with pose.



(a)



(b)



(c)

Percentage of first-match correct performance with the tied factor analysis model, combining 21 local Gabor models. (a) FERET dataset; (b) XM2VTS dataset; (c) PIE dataset.

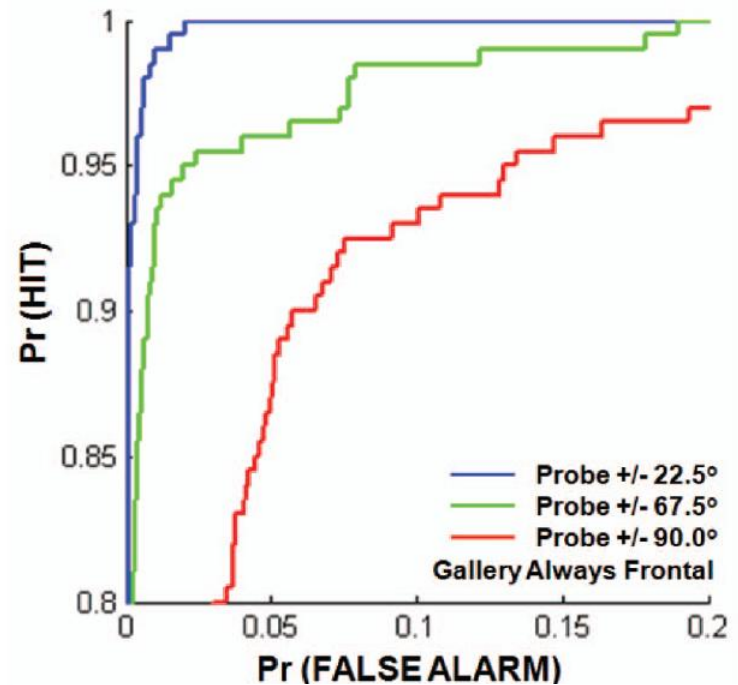
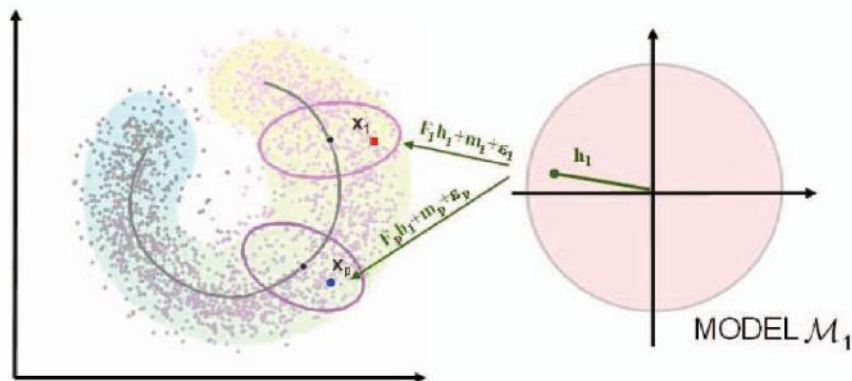
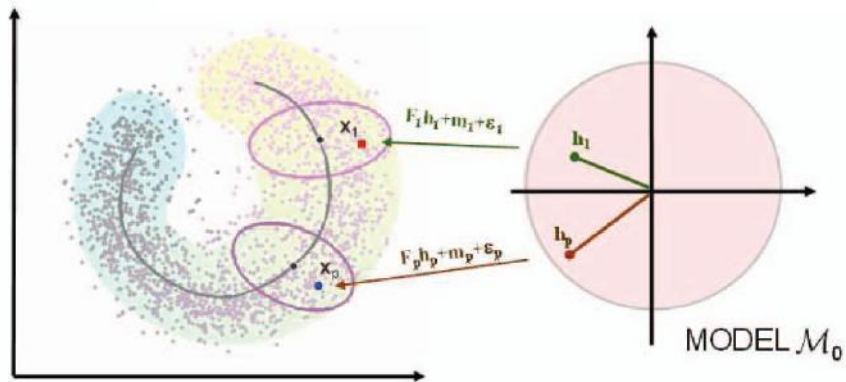
EXPERIMENT 3: FACE VERIFICATION

$$p(x_1, x_p | \mathcal{M}_0, \theta) = \int p(h_1)p(x_1 | h_1, \theta_{q_1}) dh_1 \cdot \int p(h_p)p(x_p | h_p, \theta_{q_p}) dh_p$$

$$= \mathcal{N}(x_1; m_{q_1}, F_{q_1} F_{q_1}^T + \Sigma_{q_1}) \cdot \mathcal{N}(x_p; m_{q_p}, F_{q_p} F_{q_p}^T + \Sigma_{q_p})$$

$$p(x_1, x_p | \mathcal{M}_1, \theta) = \int p(h_1)p(x_1, x_p | h_1, \theta_{q_1}, \theta_p) dh_1$$

$$= \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_p \end{bmatrix}; \begin{bmatrix} m_{q_1} \\ m_{q_p} \end{bmatrix}, \begin{bmatrix} F_{q_1} \\ F_{q_p} \end{bmatrix} \begin{bmatrix} F_{q_1} \\ F_{q_p} \end{bmatrix}^T + \begin{bmatrix} \Sigma_{q_1} & \\ & \Sigma_{q_p} \end{bmatrix}\right)$$

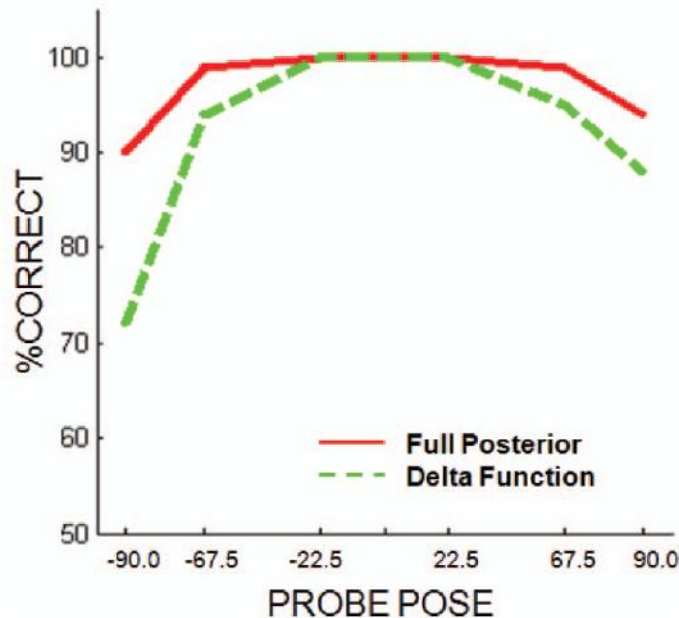


ROC curves of face verification using 21 local models.

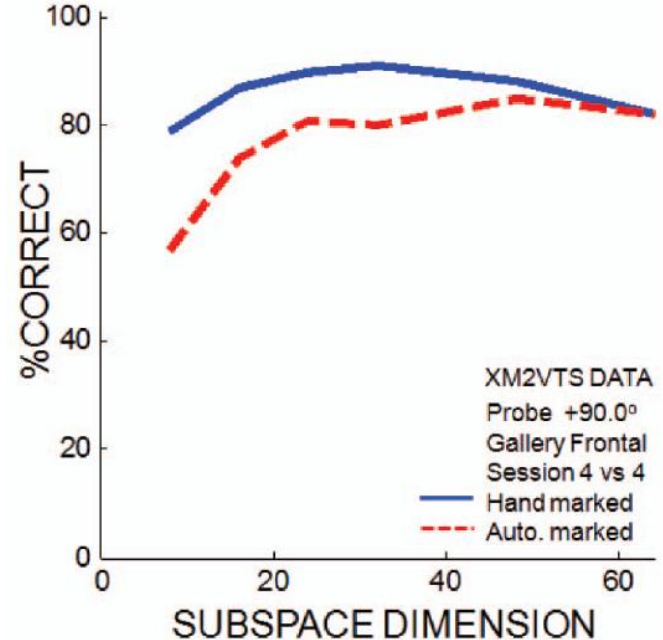
EXPERIMENTS 4 AND 5: APPROXIMATION OF EVIDENCE TERM AND AUTOMATED VERSUS MANUAL KEYPOINT DETECTION

$$p(x_1, \dots, x_N, x_p | \mathcal{M}_n, \theta) \approx p(\hat{h}_1)p(x_1 | \hat{h}_1, \theta_{q_1}) \cdots p(\hat{h}_n)p(x_n, x_p | \hat{h}_n, \theta_{q_n}, \theta_{q_p}) \cdots p(\hat{h}_N)p(x_N | \hat{h}_N, \theta_{q_N}) \\ \propto p(\hat{h}_n)p(x_p, \hat{h}_n | \theta_{q_p})$$

$$p(x_n, x_p, \hat{h}_n | \theta_{q_n}, \theta_{q_p}) = p(x_n, \hat{h}_n | \theta_{q_n}) \cdot p(x_p, \hat{h}_n | \theta_{q_p})$$



Plot of the percentage of first-match correct performance for both full and approximate (delta function) models.



Plot of the percentage of first-match correct performance for both automated and manual keypoint registration.

EXPERIMENT 6: COMPARISON TO OTHER STUDIES

Comparison of Face Identification Studies across Poses

STUDY	DATABASE	POSE DIFF ($^{\circ}$)	% CORRECT
Gross et al. [12]	FERET (100)	30 ¹	75
Gross et al. [12]	PIE (34)	6 / 62	39 / 93
Blanz et al. [5]	FRVT (87)	45	86
Zhang and Samaras [40]	CMU PIE (68) ²	45 / 90	92 / 55
Chai et al. [7]	CMU PIE (68)	16 / 45	99.85 / 89.7
Wallhoff [36]	Mugshot(100)	90	60
Maurer [20]	US ARL (90)	45	53
Kim and Kittler [15]	XM2VTS (125)	30	53
Kanade and Yanade [14]	CMU PIE (34)	45 / 67.5 / 90	100 / 80 / 40
Our method	FERET (100)	22.5 / 67.5 / 90	100 / 99 / 92
Our method	XM2VTS (100)	90	91
Our method	PIE (100)	16 / 62	100 / 91

STUDY	DATABASE	POSE DIFF ($^{\circ}$)	Pr(HIT)	Pr(FA)
Lucey and Chen [18] (EER)	FERET	60 / 90	0.9 / 0.84	0.1 / 0.16
Sanderson [30] (EER)	FERET	60	0.86	0.14
Our Method (EER)	FERET	22.5 / 67.5 / 90.0	0.99 / 0.96 / 0.925	0.01 / 0.04 / 0.075
Blanz et al. [5]	FRVT	45	0.79	0.01
Human Performance	XM2VTS	90.0	0.86	0.03
Our Method	FERET	22.5 / 67.5 / 90	0.99 / 0.93 / 0.80	0.01 / 0.01 / 0.03



CONCLUSION

- Fast
- Provides a posterior over the possible matches.
- Considers the case that the probe face is not in the database, without the need for choosing a threshold for the verification procedure.
- Only a single parameter: the dimension of the latent identity variables.
- Provides a clear way of incorporating multiple gallery or probe images.

