

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100

Abstract

While considerable progresses have been made on face recognition, age invariant face recognition (AIFR) still remains a major challenge in real world applications of face recognition systems. The major difficulty of AIFR arises from the fact that the facial appearance is subject to significant intra-personal changes caused by the aging process over time. In order to address this problem, we propose a novel deep face recognition framework to learn the age-invariant deep face features through a carefully designed CNN model. To the best of our knowledge, this is the first work to use the deep learning model to effectively address the age invariant face recognition problem. Extensive experiments are conducted on several public domain face aging datasets (MORPH Album2, FGNET, and CACD-VS) to demonstrate the effectiveness of the proposed model over the state-of-the-art. We also verify the excellent generalization of our new model on the famous LFW dataset.

1. Introduction

As an emerging research topic in face recognition community, age-invariant face recognition (AIFR) aims to address the face matching problem in the presence of remarkable aging variations [29]. It has many useful and practical applications, for example finding missing children or checking whether the same person has been issued multiple government documents. Despite the great advance in face recognition in the past decade, AIFR still remains a challenging problem in real world applications of face recognition systems. The major challenge of AIFR is mostly attributed to the significant changes in face appearance caused by aging process over time. Figure 1 is a typical example, in which the cross-age face images from the same person have significant intra-personal changes.

Most of the existing work on age-related face analysis mainly focus on age estimation [7, 8, 11, 12, 28, 39, 43] and aging modeling [19, 26, 36, 37]. There are very limited work directly on AIFR, which falls into two categories: generative approaches and discriminative approaches. The

Paper ID 1999

Anonymous CVPR submission



Figure 1. Cross-age face images for one of the subjects in the FGNET dataset [1]. One can see the significant intra-personal variation therein.

generative approaches [8, 19, 26] use face modeling to synthesis and render the probe face image to the same age as the gallery image before recognition. Due to the strong parametric assumptions and the complexity in modeling aging process, these methods are expensive to compute and the results are often unstable in real-world face recognition scenarios. Recently, the discriminative approaches [17, 20, 21, 24, 9, 10, 6] have drawn increasing attention. In [21], the authors extract the features with gradient orientation pyramid (GOP) and use support vector machine (SVM) for face classification. [20] combines the SIFT and local binary pattern (LBP) as feature representation and use a variation of random subspace LDA (RS-LDA) as classifier for AIFR. Some other variations of RS-LDA are also proposed [17, 24] to address the cross-age face recognition problem. More recent work include [9][4][10], which have notably improve the performance of AIFR. However, the existing methods in AIFR rely heavily on the hand-crafted feature descriptors to extract the dense features for age invariant face recognition, which may limit the performance of these methods. Designing an effective age invariant face features still remains an open problem in AIFR.

As one of the most promising feature learning tools nowadays, deep convolution neural networks (CNNs) have been successfully applied to a variety of problems in computer vision, including object detection and classification [18, 25, 31], and face recognition [38, 35, 32, 22, 42], etc. So it is desirable to use the deep learning model to address the AIFR problem. Surprisingly, there is no such work in

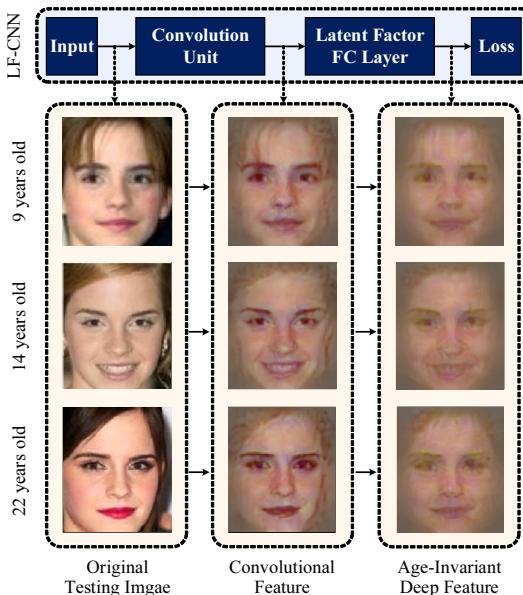


Figure 2. Cross-age faces processed by the proposed LF-CNNs. We visualize [41] the convolutional features and the age-invariant features.

the literature, to the best of our knowledge. A possible reason is the lack of a very suitable face aging dataset that can be used to train a robust deep learning model specifically for AIFR. For all the existing face aging datasets, each subject has very limited number of training samples across different ages, which are not suitable to serve as the training data in deep models. If we use the large scale web-collected face images to learn the deep CNNs, the learned deep features will inevitably contain both identity-related component (e.g. ethnicity, gender) and identity-unrelated component (e.g. age, noise) [35]. Ideally, we expect the resulting deep feature contains only the identity-related components, reducing the variations caused by the aging process as much as possible.

In this paper, we explore the use of deep CNNs in AIFR and propose a *latent factor guided CNN (LF-CNN)* framework to learn the age-invariant deep face features. Specifically, we extract the age invariant deep features from convolutional features by a carefully designed fully connected layer, termed as latent factor fully connected (LF-FC) layer. For this purpose, we develop a latent variable model, called latent identity analysis (LIA), to separate the variations caused by the aging process from the identity-related components in the convolutional features. The parameters of the LIA model are used to update the parameters of LF-FC layer. Moreover, the LIA model and the loss function in CNNs constitute the age-invariant identity loss, which is used to guide the learning of the LF-CNNs. In this way, our model is more adapted to age invariant face recognition problem, as supported by our experimental results in Section 4. Figure 2 is a visualization example of the age invariant deep features and the convolutional features, from

which we can clearly see that the convolutional features are still age sensitive, while the age-invariant features are robust to aging process.

The major contributions of this paper are summarized as follows:

- We propose a robust age invariant deep face recognition framework. To the best of our knowledge, it is the first work that utilizes the deep learning model to address the challenging AIFR task.
- Instead of directly applying deep learning model to AIFR, we propose a new model called latent factor guided convolutional neural network (LF-CNN) to specifically address the AIFR task. By coupled learning the parameters in CNNs and LIA, the age invariant deep face features can be extracted, which are more robust to the variations caused by the aging process over the time.
- Extensive experiments have shown that the propose approach significantly outperforms the state-of-the-art on all the three face aging datasets (MORPH Album2 [30], FG-NET [1] and CACD-VS [4]), even beating the human voting performance on the CACD-VS dataset. We further demonstrate the excellent generalizability of our approach on the famous LFW [15] dataset.

2. Related Work

2.1. Convolutional Networks for Face Recognition

CNNs play a significant role in recent advances of face recognition. DeepFace [38] reports a deeply-learned face representation achieves the accuracy close to human-level performance on LFW dataset [15]. [35] learns a deep CNN with the identification-verification supervisory signal and further adds supervision to early convolutional layers, greatly boosting the face recognition accuracy. FaceNet [32] achieves 99.63% verification accuracy on LFW with a deep CNN trained to directly optimize the embedding itself, rather than an intermediate bottleneck layer as in previous deep learning approaches. More recently, [22] achieves a new record in verification accuracy: 99.77% on LFW with a two-stage approach that combines a multi-patch deep CNN and deep metric learning.

2.2. Latent Variable Model

The Latent variable model finds the latent variable that are not directly observed but inferred through a statistical model from observations. Latent variable models are widely used in recommended systems [14]. A few recent work [9, 10] apply the latent variable model to face recognition

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

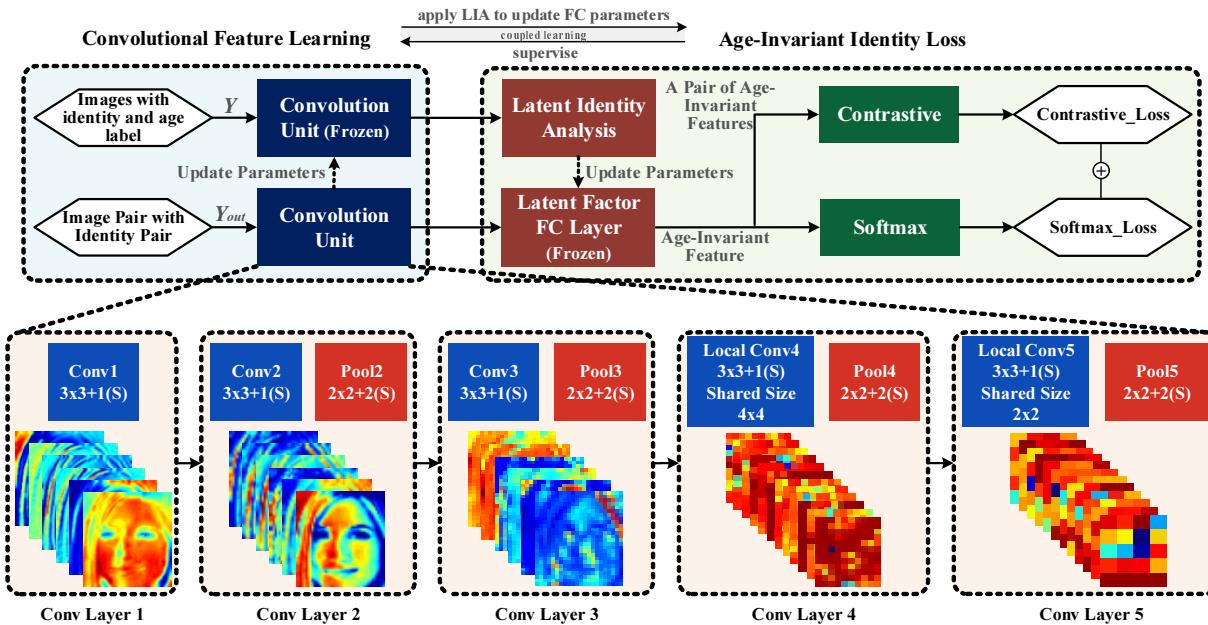


Figure 3. The architecture of the proposed LF-CNNs and its training process. *Frozen* means that the layer(s) do not update the parameters by SGD.

and achieves impressive performance. However, how to smoothly apply latent variable model in CNNs to achieve robust recognition still remains to be explored.

3. The Proposed Method

3.1. The LF-CNNs Model

The LF-CNN model is composed of two key components: convolution unit for convolution feature learning and LF-FC layer for age invariant deep feature learning. The architecture of the LF-CNNs is shown in Figure 3.

The structure of the convolution unit in LF-CNNs follows typical CNNs, alternatively stacking convolution layer, nonlinearity layer, and optional pooling layer. We construct the convolution unit with 5 convolution layers as shown in Figure 3. The convolution kernel size and stride are set as 3×3 and 1 respectively, to capture more facial details in raw images [33]. In the 4th and 5th layer, the weights of convolution are locally shared to learn different mid-level and high-level features from different regions [38]. The five convolution layers output 128, 128, 128, 256 and 256 feature maps respectively. The nonlinear function is the Parametric Rectified Linear Unit (PReLU) [13], which improves model fitting. The max pooling is used for enhancing the robustness to potential translation and subsampling.

Next we focus on the construction of the LF-FC layer. Notice that the FC layer is equivalent to the matrix multiplication: $\mathbf{F}_{fc} = \mathbf{W}\mathbf{F}_{conv} + \mathbf{b}$ where \mathbf{F}_{fc} is the output of FC layer, \mathbf{F}_{conv} is the convolutional feature, and \mathbf{W}, \mathbf{b} are the parameters of the FC layer. We leverage such equivalence property to design a set of \mathbf{W}, \mathbf{b} that can extract the

age-invariant feature from \mathbf{F}_{conv} . Note that, no nonlinearity layer will cascade to the LF-FC layer according to the LIA model. The feature dimension of \mathbf{F}_{fc} is 512. Instead of iteratively updating the all the parameters in LF-CNNs by stochastic gradient descent (SGD), we design an *Latent Identity Analysis (LIA)* method to learn \mathbf{W} and \mathbf{b} for the LF-CNN model, as elaborated in the following subsection. For the parameters of the convolution unit, we fix \mathbf{W}, \mathbf{b} to update them via standard SGD.

3.2. Latent Identity Analysis

The latent identity analysis model can infer the unobserved latent factors (one of them being the identity factor) from the observed features in a supervised fashion. The general model of latent identity analysis is formulated as

$$\mathbf{v} = \sum_{i=1}^d \mathbf{U}_i \mathbf{x}_i + \bar{\mathbf{v}} \quad (1)$$

where $\mathbf{v} \in \mathbb{R}^{n \times 1}$ denotes the observed facial features, $\mathbf{U}_i \in \mathbb{R}^{n \times p_i}$ is the corresponding matrix whose columns span the subspace of different variation and $\mathbf{x}_i \in \mathbb{R}^{p_i \times 1}$ denotes the latent variable with prior zero-mean Gaussian distribution. $\bar{\mathbf{v}} \in \mathbb{R}^{n \times 1}$ is the mean of all the facial features. The intuition behind this model is very clear. Each facial feature is viewed as the combination of different components according to different supervised signals. Such an idea is useful to achieve robust face recognition in practical. In the cross-age face recognition, we usually decompose the facial features into two latent components and a noise variable. So the model can be simplified as

324 $\mathbf{v} = \mathbf{U}_{\text{id}}\mathbf{x}_{\text{id}} + \mathbf{U}_{\text{ag}}\mathbf{x}_{\text{ag}} + \mathbf{U}_{\text{e}}\mathbf{x}_{\text{e}} + \bar{\mathbf{v}}$ where $\mathbf{x}_{\text{id}}, \mathbf{x}_{\text{ag}}$ satisfy standard Gaussian distribution $\mathcal{N}(0, I)$ and \mathbf{x}_{e} satisfies $\mathcal{N}(0, \sigma^2 I)$. Note that \mathbf{U}_{e} is set to be a unix matrix since \mathbf{x}_{e} stands for noise. $\mathbf{U}_{\text{id}}\mathbf{x}_{\text{id}}$ is the identity-related component which is key to achieve age invariant face recognition. $\mathbf{U}_{\text{ag}}\mathbf{x}_{\text{ag}}$ is the age-related component, representing the age variations.

331 A set of model parameters $\theta = \{\mathbf{U}_{\text{id}}, \mathbf{U}_{\text{ag}}, \sigma^2, \bar{\mathbf{v}}\}$ can be learned by maximizing the following maximum likelihood function $\mathcal{L}(\theta)$:

$$\begin{aligned}\mathcal{L}(\theta) &= \sum_{i,j} \mathcal{L}_{i,j}(\theta) \\ &= \sum_{i,j} \ln P(\mathbf{v}_i^j, \mathbf{x}_{\text{id},j}, \mathbf{x}_{\text{ag},j} | \theta)\end{aligned}\quad (2)$$

340 where \mathbf{v}_i^j is the feature of the i th subject at j th age group. **341** $\mathbf{x}_{\text{id},i}$ and $\mathbf{x}_{\text{ag},j}$ are the corresponding identity and age factors respectively. The summation is over all the available **342** training samples from all subjects at all age groups. We use the **343** expectation-maximization (EM) algorithm to estimate these **344** model parameters. To perform the EM algorithm, the **345** Q function is given by

$$\begin{aligned}\mathcal{Q}(\theta, \theta^{(i)}) &= \sum_{i,j} E\{\ln P(\mathbf{v}_i^j, \mathbf{x}_{\text{id},j}, \mathbf{x}_{\text{ag},j} | \theta)\} \\ &= \sum_{i,j} \int_{\mathbf{x}_{\text{id},i}, \mathbf{x}_{\text{ag},j}} P(\mathbf{x}_{\text{id},j}, \mathbf{x}_{\text{ag},j} | \theta^{(i)}, \mathbf{V}) \mathcal{L}_{i,j}(\theta)\end{aligned}\quad (3)$$

352 where \mathbf{V} denotes the observed features of all the training **353** samples, $\theta^{(i)}$ is the given model parameters and θ is the **354** parameters to be estimated. With the given parameter $\theta^{(i)}$, we **355** can compute the posterior distribution of the latent variables **356** $P(\mathbf{x}_{\text{id},j}, \mathbf{x}_{\text{ag},j} | \theta^{(i)}, \mathbf{V})$. With the given posterior **357** distribution, we can maximize the Q function to obtain a new **358** θ . The EM algorithm is performed in an iterative fashion.

359 **E Step.** Given the model parameter $\theta^{(i)}$ and training **360** data $\mathbf{V} = \{\mathbf{v}_i^j\}_{i=1, \dots, N; j=1, \dots, M}$, we first compute **361** all the necessary first and second conditional moments of **362** $P(\mathbf{x}_{\text{id},j} | \theta^{(i)}, \mathbf{V})$ and $P(\mathbf{x}_{\text{ag},j} | \theta^{(i)}, \mathbf{V})$ for the **363** posterior distribution $P(\mathbf{x}_{\text{id},j}, \mathbf{x}_{\text{ag},j} | \theta^{(i)}, \mathbf{V})$:

$$\mu_1(\mathbf{x}_{\text{id},i}) = \frac{\mathbf{U}_{\text{id}}\Sigma^{-1}}{N_i} \sum_{k=1}^{N_i} (\mathbf{v}_i^k - \bar{\mathbf{v}}) \quad (4)$$

$$\mu_1(\mathbf{x}_{\text{ag},j}) = \frac{\mathbf{U}_{\text{ag}}\Sigma^{-1}}{M_j} \sum_{k=1}^{M_j} (\mathbf{v}_k^j - \bar{\mathbf{v}}) \quad (5)$$

$$\mu_2(\mathbf{x}_{\text{id},i}, \mathbf{x}_{\text{id},i}) = \frac{\mathbf{I} - \mathbf{U}_{\text{id}}^T \Sigma^{-1} \mathbf{U}_{\text{id}}}{N_i} + \mu_1(\mathbf{x}_{\text{id},i})(\mu_1(\mathbf{x}_{\text{id},i}))^T \quad (6)$$

$$\mu_2(\mathbf{x}_{\text{ag},j}, \mathbf{x}_{\text{ag},j}) = \frac{\mathbf{I} - \mathbf{U}_{\text{ag}}^T \Sigma^{-1} \mathbf{U}_{\text{ag}}}{M_j} + \mu_1(\mathbf{x}_{\text{ag},j})(\mu_1(\mathbf{x}_{\text{ag},j}))^T \quad (7)$$

$$\mu_2(\mathbf{x}_{\text{id},i}, \mathbf{x}_{\text{ag},j}) = -\frac{\mathbf{U}_{\text{ag}}^T \Sigma^{-1} \mathbf{U}_{\text{id}}}{\sqrt{N_i M_j}} + \mu_1(\mathbf{x}_{\text{id},i})(\mu_1(\mathbf{x}_{\text{ag},j}))^T \quad (8)$$

$$\mu_2(\mathbf{x}_{\text{ag},j}, \mathbf{x}_{\text{id},i}) = -\frac{\mathbf{U}_{\text{id}}^T \Sigma^{-1} \mathbf{U}_{\text{ag}}}{\sqrt{N_i M_j}} + \mu_1(\mathbf{x}_{\text{ag},j})(\mu_1(\mathbf{x}_{\text{id},i}))^T \quad (9)$$

where $\Sigma = \sigma^2 \mathbf{I} + \mathbf{U}_{\text{id}}\mathbf{U}_{\text{id}}^T + \mathbf{U}_{\text{ag}}\mathbf{U}_{\text{ag}}^T$, N_i and M_j are the numbers of training samples for the i th subject and the k th age group, respectively.

M Step. We maximize the Q function to estimate the $\theta^{(i+1)}$. The maximization is shown as follows:

$$\theta^{(i+1)} = \arg \max_{\theta} \mathcal{Q}(\theta, \theta^{(i)}) \quad (10)$$

To solve this optimization, the model parameter $\theta^{(i+1)}$ is given by

$$\begin{aligned}\mathbf{U}_{\text{id}} &= (\mathbf{C} - \mathbf{D}\mathbf{B}^{-1}\mathbf{E})(\mathbf{A} - \mathbf{F}\mathbf{B}^{-1}\mathbf{E})^{-1} \\ \mathbf{U}_{\text{ag}} &= (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{F})(\mathbf{B} - \mathbf{E}\mathbf{A}^{-1}\mathbf{F})^{-1} \\ \sigma^2 &= \frac{1}{Nn} \sum_{i,j} \{(\mathbf{v}_i^j - \bar{\mathbf{v}} - \mathbf{U}_{\text{id}}\mu_1(\mathbf{x}_{\text{id},i}) \\ &\quad - \mathbf{U}_{\text{ag}}\mu_1(\mathbf{x}_{\text{ag},j}))^T(\mathbf{v}_i^j - \bar{\mathbf{v}})\}\end{aligned}\quad (11)$$

in which

$$\begin{aligned}\mathbf{A} &= \sum_{ij} \mu_2(\mathbf{x}_{\text{id},i}, \mathbf{x}_{\text{id},i}), \mathbf{B} = \sum_{ij} \mu_2(\mathbf{x}_{\text{ag},j}, \mathbf{x}_{\text{ag},j}), \\ \mathbf{C} &= \sum_{ij} (\mathbf{v}_i^j - \bar{\mathbf{v}})(\mu_1(\mathbf{x}_{\text{id},i}))^T, \mathbf{D} = \sum_{ij} (\mathbf{v}_i^j - \bar{\mathbf{v}})(\mu_1(\mathbf{x}_{\text{ag},j}))^T, \\ \mathbf{E} &= \sum_{ij} \mu_2(\mathbf{x}_{\text{ag},j}, \mathbf{x}_{\text{id},i}), \mathbf{F} = \sum_{ij} \mu_2(\mathbf{x}_{\text{id},i}, \mathbf{x}_{\text{ag},j}).\end{aligned}\quad (12)$$

After the model parameters and the posterior distribution are estimated by the EM algorithm, the identity factor for the i th subject can be inferred by the first moment of $\mathbf{x}_{\text{id},i}$, namely $\mu_1(\mathbf{x}_{\text{id},i})$. One can observe that the form of $\mu_1(\mathbf{x}_{\text{id},i})$ is actually identical to an LF-FC layer with the following parameters:

$$\mathbf{W} = \mathbf{U}_{\text{id}}^T \Sigma^{-1}, \mathbf{b} = -\mathbf{U}_{\text{id}}^T \Sigma^{-1} \bar{\mathbf{v}} \quad (13)$$

where $\bar{\mathbf{v}} = \frac{1}{N} \sum_{i,j} \mathbf{v}_i^j$ and $\Sigma = \sigma^2 \mathbf{I} + \mathbf{U}_{\text{id}}\mathbf{U}_{\text{id}}^T + \mathbf{U}_{\text{ag}}\mathbf{U}_{\text{ag}}^T$. The feature \mathbf{v} comes from the convolution unit. Because \mathbf{v} is updated with new parameters of the identity factor guided FC layer, \mathbf{W}, \mathbf{b} are also updated using the new \mathbf{v} . The procedure iteratively goes on. The detailed learning framework is elaborated in the next subsection.

3.3. Learning Framework

In LF-CNN model, the convolution unit maps a raw input image \mathbf{F}_I to convolutional feature \mathbf{F}_{conv} by $\mathbf{F}_{\text{conv}} = f(\mathbf{F}_I)$, and then the LF-FC layer computes the age-invariant feature \mathbf{F}_{fc} via $\mathbf{F}_{\text{fc}} = g(\mathbf{F}_{\text{conv}})$. The age-invariant feature \mathbf{F}_{fc} is used for the AIFR. The parameters in the convolution unit and the LF-FC layer are denoted by θ_f, θ_g respectively. In our framework, $f(\cdot)$ and $g(\cdot)$ characterize different properties of AIFR and the original learning methods based on CNNs are not appropriate, we adopt a

Algorithm 1 Coupled Learning Algorithm for LF-CNNs

Input: Outside training data \mathbf{Y}_{out} with identity label, cross-age training data \mathbf{Y} with both age and identity label.
Output: The parameters θ_f and θ_g .

- 1: $i \leftarrow 0$.
- 2: Initialize the parameters $\theta_f^{(1)}$ by Xaiver filter.
- 3: Initialize $\theta_g^{(1)} = \{\mathbf{U}_{\text{id}}, \mathbf{U}_{\text{ag}}, \sigma^2, \bar{v}\}$ where $\sigma^2 = 0.1$, $\bar{v} = \mathbf{0}$ and \mathbf{U}_{id} , \mathbf{U}_{ag} are randomly initialized from -0.1 to 0.1 .
- 4: Compute $\mathbf{W}^{(1)}, \mathbf{b}^{(1)}$ for the LF-FC layer via Eq. (13).
- 5: **while** not converge **do**
- 6: $i \leftarrow i + 1$.
- 7: Fix the $\theta_g^{(i)}$, and train the LF-CNNs with the outside data \mathbf{Y}_{out} to update $\theta_f^{(i+1)}$ from $\theta_f^{(i)}$.
- 8: Fix the $\theta_f^{(i+1)}$, and input the training data \mathbf{Y} to obtain the convolutional features \mathbf{F}_{conv} , which are taken as the observed features \mathbf{V} .
- 9: Update the parameters $\theta_g^{(i+1)}$ from $\theta_g^{(i)}$ via Eq. (3).
- 10: Compute the parameters $\mathbf{W}^{(i+1)}, \mathbf{b}^{(i+1)}$ for the LF-FC layer via the Eq. (13).
- 11: **end while**

coupled learning framework to optimize the LF-CNN model. Concretely, after initializing the networks, the parameters θ_f of convolutional unit is updated by SGD with the fixed θ_g . Then we fix θ_f and use the LIA model to learn the parameters θ_g for the LF-FC layer. We alternatively update θ_f and θ_g until the stopping condition is satisfied. The procedure is summarized in Algorithm 1.

Learning parameters for LF-FC layer. We use the training data \mathbf{Y} with both age and identity label to train the LF-FC layer. Specifically, the convolutional feature \mathbf{F}_{fc} is taken as the observed features \mathbf{V} in LIA. The LIA model learns the parameters $\theta_g = \{\mathbf{U}_{\text{id}}, \mathbf{U}_{\text{ag}}, \sigma^2, \bar{v}\}$, and then computes \mathbf{W}, \mathbf{b} for the LF-FC layer via Eq. (13).

Learning parameters for convolution unit. To learn the parameters θ_f for the convolution unit, we need to fix the parameters θ_g of the LF-FC layer. Then we use the SGD to train the LF-CNNs with outside training data \mathbf{Y}_{out} . Note that, we use both the softmax loss and the contrastive loss to strength the supervision in learning, similar to [34].

3.4. Discussion

The learning process of LF-CNNs has the following advantages. First, the coupled learning is very beneficial to AIFR. The joint objective functions consists of minimizing the classification error (softmax loss and contrastive loss) and maximizing the likelihood probability that the training samples are generated by the latent factors. The former aims to learn discriminant feature representations for classification while the later improves the robustness of age-invariant features. Both of them consistently contribute to the gain in the AIFR task. Second, we update the LF-FC layer by LIA instead of SGD, which largely reduces the parameter scale

and prevents potential overfitting. LIA plays an essential role in LF-CNNs by inferring the effective identity factor to guide the parameter estimation of the LF-FC layer.

4. Experiments

In this section, we conduct extensive experiments to demonstrate the effectiveness of the proposed LF-CNNs on several challenging face aging databases, including MORPH Album 2 [30] (the largest face aging database available in the public domain), FG-NET [1] (a public-domain face aging dataset), and the subset of Cross-Age Celebrity Dataset (CACD-VS) [4]. To further demonstrate the generalization of our model, we also evaluate our model on the Labeled Faces in the Wild (LFW) database [15].

4.1. Implementation Details

Preprocessing. For each face image, we use the recently proposed algorithms [2, 44, 3] to detect the facial landmarks in images. Then the face are globally cropped to 112×96 according to the 5 facial landmarks (two eyes, nose and two mouth corners) by similarity transformation.

Training Data. The training data used in this paper is composed of two parts: outside training data \mathbf{Y}_{out} (that only contains the identity information) and face aging training data \mathbf{Y} (that contains both the age and identity information). For the outside training data \mathbf{Y}_{out} , we use the large scale web-collected face data, including CASIA-WebFace [40], CACD [4], and Celebrity+ [23]). When testing our model on CACD-VS, we remove all the identities in CACD from the training data. The images are horizontally flipped for data augmentation. For the face aging training data \mathbf{Y} , we use the MORPH Album 2 dataset, as described in Section 4.2.

Detailed setting in LF-CNN model. We implement the LF-CNN model using the Caffe library [16] with our modifications. Unless otherwise specified, the batch size is 150. When training convolution unit with outside data, the learning rate is 1e-1, 1e-2, 1e-3 and is switched when the error plateaus. The total number of epochs is about 12 for our model.

Classifier. To better evaluate the performance of the proposed age invariant deep face features, our model uses the simple Euclidean Distance and the Nearest Neighbor rule as the classifier.

4.2. Experiments on the MORPH Album 2 Dataset

The MORPH Album 2 database is the largest face aging dataset available in the public domain, consisting of about 78,000 face images of 20,000 persons with age ranging from 16 to 77. Following the same training and testing split scheme in [9], 10,000 subjects are used for training and the remaining 10,000 subjects are used for testing. There

is no overlapping subject between the training set and the testing set. For each subject, two face images with the youngest age and the oldest age are selected as gallery and probe set respectively. For fair comparison, we also train a baseline CNN model with the same networks as LF-CNNs, and learn all the parameters by SGD. The experimental results are shown in Table 1.

Method	Rank-1 Identification Rates
HFA (2013) [9]	91.14%
CARC (2014) [4]	92.80%
MEFA (2015) [10]	93.80%
MEFA+SIFT+MLBP (2015) [10]	94.59%
CNN-baseline	89.68%
CNN-baseline (fine-tuned by MORPH training data)	95.13%
LF-CNNs (fine-tuned by MORPH training data)	97.51%

Table 1. Performance of different methods on MORPH.

In Table 1, we compare our LF-CNN model against (i) the CNN-baseline model, (ii) the CNN-baseline model (finetuned by MORPH training data), and (iii) several recently developed top-performing AIFR algorithms in the literature. From these results, we have the following observations. First, the result of the CNN-baseline is only 89.68%, which is inferior to the other results in Table 1. This confirms that directly applying the deep CNN model to address the AIFR problem is indeed not a good choice. Second, the performance of CNN-baseline can be improve to 95.13% by fine-tuning with the additional MORPH training data. However, this result (95.13%) is near to the top-performing result in the literature (94.59%). The lack of a significant improvement over the state-of-the-art reflects the limitation of the CNN-baseline model. So it is desirable to design a new deep CNN model to address the AIFR problem. Finally, it is encouraging to see that the proposed LF-CNN model obtains a significant performance improvement over the other results in Table 1, demonstrating a new state-of-the-art (97.51%) on the MORPH Album 2 database.

To further evaluate the performance of our LF-CNN model, we design an experiment to report the recognition results of our model in each iteration of the coupled learning process, as illustrated in Figure 4. The two parameters θ_f and θ_g are updated in a step-wise manner. With one fixed, update the other one, and vice versa. Figure 4 clearly shows that the coupled learning process consistently contributes to the performance improvement of AIFR, converging to a good result quickly.

Figure 5 shows some examples of the failed retrievals using our approach in MORPH Album 2 dataset. We can see that although our results are incorrect in these cases, the probe images appear to be more similar to the incorrect retrievals than the gallery images.

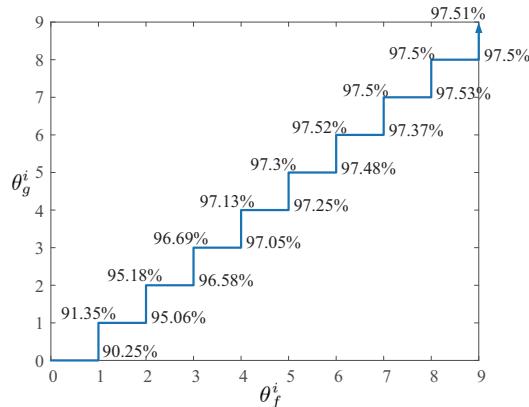


Figure 4. The recognition rates in each iteration of the coupled learning algorithm.



Figure 5. Some failed retrievals in MORPH Album 2. The first row: the probe images. The second row: the incorrect retrievals using our approach. The third row: the corresponding gallery images for the probe images.

4.3. Experiments on the FG-NET Dataset

The FG-NET dataset consists of 1002 face images from 82 different subjects, with each subject having multiple face images at different ages (ranging from 0 to 69). Following the testing scheme in [20], we compare our LF-CNNs with the state-of-the-art approaches on this dataset. The comparative results are reported in Table 2.

Method	Rank-1 Identification Rates
Park et al. (2010) [9]	37.4%
Li et al. (2011) [4]	47.5%
HFA (2013) [9]	69.0%
MEFA (2015) [10]	76.2%
CNN-baseline	84.4%
LF-CNNs	88.1%

Table 2. Performance of different methods on FG-NET.

As can be seen from Table 2, LF-CNNs achieve the highest recognition accuracy (88.1%) among all the results, significantly outperforming the top-performing result (76.2%) in [10] by 11.9%. Moreover, the proposed LF-CNNs outperform the CNN-baseline method by a clear margin. This

648 confirms what we observe from the MORPH Album 2
 649 dataset. Interestingly, MORPH and FG-NET have different
 650 age distributions. In FGNET, roughly 61% samples are less
 651 than 16 years old. But for the MORPH dataset, all the per-
 652 sons are more than 16 years old. So it is desirable to exploit
 653 the influence of different age distributions on the proposed
 654 approach. In Table 3 we give the rank-1 identification rates
 655 in different age groups.
 656

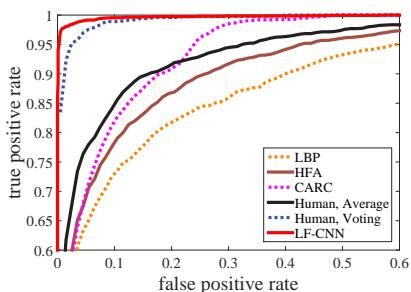
Age group	Amount	CNN-baseline	LF-CNNs
0 - 4	193	51.81%	60.10%
5 - 10	218	84.86%	88.53%
11 - 16	201	91.04%	94.03%
17 - 24	182	94.51%	97.80%
25 - 69	208	99.04%	99.52%
0 - 16	612	76.47%	81.37%
17 - 69	390	96.93%	98.72%

664 Table 3. Performance of different age groups on FG-NET.
 665

666 The results in Table 3 show that the proposed LF-CNNs
 667 consistently outperforms the CNN-baseline model on all the
 668 age groups. This further confirms the advantage of our LF-
 669 CNN model over the CNN-baseline model in AIFR task.
 670

671 4.4. Experiments on the CACD Verification Subset

672 The CACD dataset is a recently released dataset for
 673 AIFR, containing 163,446 images from 2,000 celebrities
 674 with labeled ages. It includes varying illumination, pose
 675 variation and makeup and better simulates practical scenari-
 676 o. However, the entire CACD dataset contains some incor-
 677 rectly labeled samples, and some duplicate images. Follow-
 678 ing the state-of-the-art configuration [4], we test LF-CNNs
 679 on a subset of CACD [4], CACD-VS, which consists of
 680 4000 image pairs (2000 positive pairs and 2000 negative
 681 pairs) and have been carefully annotated. We follow the
 682 same training strategy as in section 4.3. Note that, the iden-
 683 tities in CACD-VS are excluded from the outside training
 684 data in this experiment. So only 400,000 training samples
 685 are used. According to the ten-fold cross-validation rule, we
 686 calculate the Euclidean Distance of each pairs and choose
 687 the best threshold using nine training folds, then testing on
 688 the leftover fold. We compute the face verification rate and
 689 compare our result with the existing methods in this dataset,
 690 as shown in Figure 6 and Table 4.
 691



692 Figure 6. ROC comparisons of different methods on CACD-VS.
 693

694 Again, the proposed LF-CNN model significantly out-
 695 performs all the published results in this dataset, even sur-
 696 passing the human-level performance with a clear margin.
 697 It further demonstrates the effectiveness of the proposed age
 698 invariant deep features.
 699

Method	Acc.
High-Dimensional LBP [5]	81.6%
HFA (2013) [9]	84.4%
CARC (2014) [4]	87.6%
Human, Average (2013)	85.7%
Human, Voting (2015)	94.2%
LF-CNNs	98.5%

700 Table 4. Performance of different methods on CACD-VS.
 701

702 4.5. Experiments on the LFW Dataset

703 To evaluate the generalization performance of LF-CNNs,
 704 we further conduct an experiment on the famous LFW
 705 dataset [15]. This dataset contains 13,233 face images from
 706 5749 different subjects, collecting from uncontrolled condi-
 707 tions. Following the unrestricted with labeled outside data
 708 protocol [15], we train on the outside dataset and test on
 709 6,000 face pairs. People overlapping between the outside
 710 training data and the LFW testing data are excluded. We re-
 711 spectively train 25 networks with 25 different image patch-
 712 es, and concatenate the output features from these networks
 713 into a long feature vector. We then apply PCA on the long
 714 feature vector to obtain a compact feature vector for classifi-
 715 cation. In Table 5 we compare our results against the recent
 716 state-of-the-art results. From the results we can see that our
 717 approach can obtain comparable results to the state-of-the-
 718 art approaches using relatively small training data, demon-
 719 strating the excellent generalization ability of our approach.
 720

Method	Images	Networks	Acc.
DeepFace [38]	4M	3	97.35%
DeepID-2+ [35]	-	25	99.47%
FaceNet [32]	200M	1	99.65%
Deep Embedding [22]	1.2M	10	99.77%
Deep FR [27]	2M	1	98.95%
LF-CNNs(single model)	700K	1	99.10%
LF-CNNs(ensemble model)	700K	25	99.50%

737 Table 5. Performance of different methods on LFW.
 738

744 5. Conclusions

745 In this paper, we have proposed an age invariant deep
 746 face recognition framework, referred to as LF-CNNs. Un-
 747 like the existing deep learning models in face recogni-
 748 tion community, the proposed new model constructs a latent i-
 749 dentity analysis (LIA) module to guide the learning of the
 750 CNNs parameters. By coupled learning the CNNs parame-
 751 ters and the LIA parameters, our model can extract the age
 752

756 invariant deep face features, which are well suitable for the
 757 AIFR task. Extensive experiments are conducted on several
 758 public-domain face aging databases to demonstrate the sig-
 759 nificant performance improvement of our new model over
 760 the state-of-the-art. We have also performed experiments
 761 on the famous LFW dataset to demonstrate the excellent
 762 generalization ability of our new model.
 763

References

- 756 [1] Fg-net aging database. In <http://www.fgnet.rsunit.com/>, 2010. 1, 2, 5
- 757 [2] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Robust discriminative response map fitting with constrained local models. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3444–3451. IEEE, 2013. 5
- 758 [3] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Incremental face alignment in the wild. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1859–1866. IEEE, 2014. 5
- 759 [4] B.-C. Chen, C.-S. Chen, and W. Hsu. Face recognition and retrieval using cross-age reference coding with cross-age celebrity dataset. *IEEE TMM*, 17(6):804–815, 2015. 1, 2, 5, 6, 7
- 760 [5] D. Chen, X. Cao, F. Wen, and J. Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3025–3032. IEEE, 2013. 7
- 761 [6] L. Du and H. Ling. Cross-age face verification by coordinating with cross-face age verification. In *CVPR*, 2015. 1
- 762 [7] Y. Fu and T. S. Huang. Human age estimation with regression on discriminative aging manifold. *IEEE TMM*, 10(4):578–584, 2008. 1
- 763 [8] X. Geng, Z.-H. Zhou, and K. Smith-Miles. Automatic age estimation based on facial aging patterns. *IEEE TPAMI*, 29(12):2234–2240, 2007. 1
- 764 [9] D. Gong, Z. Li, D. Lin, J. Liu, and X. Tang. Hidden factor analysis for age invariant face recognition. In *ICCV*, 2013. 1, 2, 5, 6, 7
- 765 [10] D. Gong, Z. Li, D. Tao, J. Liu, and X. Li. A maximum entropy feature descriptor for age invariant face recognition. In *CVPR*, 2015. 1, 2, 6
- 766 [11] G. Guo, Y. Fu, C. R. Dyer, and T. S. Huang. Image-based human age estimation by manifold learning and locally adjusted robust regression. *IEEE TIP*, 17(7):1178–1188, 2008. 1
- 767 [12] G. Guo, G. Mu, Y. Fu, and T. S. Huang. Human age estimation using bio-inspired features. In *CVPR*, 2009. 1
- 768 [13] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *arXiv preprint arXiv:1502.01852*, 2015. 3
- 769 [14] T. Hofmann and J. Puzicha. Latent class models for collaborative filtering. In *IJCAI*, 1999. 2
- 770 [15] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007. 2, 5, 7
- 771 [16] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 5
- 772 [17] B. Klare and A. K. Jain. Face recognition across time lapse: On learning feature subspaces. In *IJCB*, 2011. 1
- 773 [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1
- 774 [19] A. Lanitis, C. J. Taylor, and T. F. Cootes. Toward automatic simulation of aging effects on face images. *IEEE TPAMI*, 24(4):442–455, 2002. 1
- 775 [20] Z. Li, U. Park, and A. K. Jain. A discriminative model for age invariant face recognition. *Information Forensics and Security, IEEE Transactions on*, 6(3):1028–1037, 2011. 1, 6
- 776 [21] H. Ling, S. Soatto, N. Ramanathan, and D. W. Jacobs. Face verification across age progression using discriminative methods. *Information Forensics and Security, IEEE Transactions on*, 5(1):82–91, 2010. 1
- 777 [22] J. Liu, Y. Deng, B. Tao, W. Zhengping, and C. Huang. Targeting ultimate accuracy: Face recognition via deep embedding. *arXiv preprint arXiv:1506.07310*, 2015. 1, 2, 7
- 778 [23] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. *arXiv preprint arXiv:1411.7766*, 2014. 5
- 779 [24] C. Otto, H. Han, and A. Jain. How does aging affect facial components? In *ECCV*, 2012. 1
- 780 [25] W. Ouyang, X. Wang, X. Zeng, S. Qiu, P. Luo, Y. Tian, H. Li, S. Yang, Z. Wang, C.-C. Loy, et al. Deepid-net: Deformable deep convolutional neural networks for object detection. In *CVPR*, 2015. 1
- 781 [26] U. Park, Y. Tong, and A. K. Jain. Age-invariant face recognition. *IEEE TPAMI*, 32(5):947–954, 2010. 1
- 782 [27] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. *Proceedings of the British Machine Vision*, 2015. 7
- 783 [28] N. Ramanathan and R. Chellappa. Face verification across age progression. *IEEE TIP*, 15(11):3349–3361, 2006. 1
- 784 [29] N. Ramanathan, R. Chellappa, and S. Biswas. Computational methods for modeling facial aging: A survey. *Journal of Visual Languages & Computing*, 20(3):131–144, 2009. 1
- 785 [30] K. Ricanek Jr and T. Tesafaye. Morph: A longitudinal image database of normal adult age-progression. In *FG*, 2006. 2, 5
- 786 [31] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, pages 1–42, 2014. 1
- 787 [32] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 1, 2, 7
- 788 [33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3
- 789 [34] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems*, pages 1988–1996, 2014. 5
- 790 [35] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. In *CVPR*, 2015. 1, 2, 7
- 791 [36] J. Suo, X. Chen, S. Shan, and W. Gao. Learning long term face aging patterns from partially dense aging databases. In *ICCV*, 2009. 1
- 792 [37] J. Suo, S.-C. Zhu, S. Shan, and X. Chen. A compositional and dynamic model for face aging. *IEEE TPAMI*, 32(3):385–401, 2010. 1
- 793 [38] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, 2014. 1, 2, 3, 7
- 794 [39] S. Yan, H. Wang, X. Tang, and T. S. Huang. Learning auto-structured regressor from uncertain nonnegative labels. In *ICCV*, 2007. 1
- 795 [40] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 5
- 796 [41] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014*, pages 818–833. Springer, 2014. 2
- 797 [42] E. Zhou, Z. Cao, and Q. Yin. Naive-deep face recognition: Touching the limit of lfw benchmark or not? *arXiv preprint arXiv:1501.04690*, 2015. 1
- 798 [43] S. K. Zhou, B. Georgescu, X. S. Zhou, and D. Comaniciu. Image based regression using boosting method. In *CVPR*, 2005. 1
- 799 [44] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2879–2886. IEEE, 2012. 5
- 800 [810] 811 [812] 813 [814] 815 [816] 817 [818] 819 [820] 821 [822] 823 [824] 825 [826] 827 [828] 829 [830] 831 [832] 833 [834] 835 [836] 837 [838] 839 [840] 841 [842] 843 [844] 845 [846] 847 [848] 849 [850] 851 [852] 853 [854] 855 [856] 857 [858] 859 [860] 861 [862] 863