

Shared Representation Learning for Heterogeneous Face Recognition

Dong Yi, Zhen Lei, Shengcai Liao and Stan Z. Li

National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China, 100190

Abstract. After intensive research, heterogenous face recognition is still a challenging problem. The main difficulties are owing to the complex relationship between heterogenous face image spaces. The heterogeneity is always tightly coupled with other variations, which makes the relationship of heterogenous face images highly nonlinear. Many excellent methods have been proposed to model the nonlinear relationship, but they apt to overfit to the training set, due to limited samples. Inspired by the unsupervised algorithms in deep learning, this paper proposes an novel framework for heterogeneous face recognition. We first extract Gabor features at some localized facial points, and then use Restricted Boltzmann Machines (RBMs) to learn a shared representation locally to remove the heterogeneity around each facial point. Finally, the shared representations of local RBMs are connected together and processed by PCA. Two problems (Sketch-Photo and NIR-VIS) and three databases are selected to evaluate the proposed method. For Sketch-Photo problem, we obtain perfect results on the CUFS database. For NIR-VIS problem, we produce new state-of-the-art performance on the CASIA HFB and NIR-VIS 2.0 databases.

Keywords: Face Recognition, Restricted Boltzmann Machines, Sketch, Near Infrared

1 Introduction

The core of heterogenous face recognition [1] is face matching across modalities. Although the original definition of heterogeneous face recognition is broad, the two hottest problems about this topic are Sketch-Photo [2] face recognition and NIR-VIS (Near Infrared-Visual) [3] face recognition. This paper will also take these two problems as examples to verify the proposed method.

Initially, heterogeneous face recognition was proposed to appeal requirements in practical applications. Sketch-Photo matching is often required in law enforcement when the photo of suspect is unavailable. NIR-VIS matching module can make VIS face recognition system work in dark environment using NIR imaging device. After several research groups were attracted to this topic, many good methods have been proposed and these methods quickly spread to other cross-modal problems, such as face hallucination [4], pedestrian detection [5] and so on.

It has been shown in existing works that the relationship of face images between different modalities is very complex, therefore nonlinear methods usually have better performance than linear methods. Taking NIR-VIS as an example, the effect of spectrum is tightly coupled with other variations of face image, such as 3D shape, pose, identity and so on, which makes the relationship of face images under different spectrums highly nonlinear and varying with respect to locations. Among existing methods, the most successful category is learning two mappings (linear or nonlinear) to project the heterogeneous face images into a common space [6][7]. Limited by the number of training samples, this kind of methods have many regularization terms, so need careful parameter tuning to achieve good performance.

From 2006 to now, unsupervised pre-training has obtained great success in deep learning [8]. One of the most popular unsupervised learning method in deep learning is Restricted Boltzmann Machine (RBM) [9], which is a generative stochastic neural network that can learn a probability distribution of input data. To improve the generalization of existing methods and make the training process easily, this paper propose a framework based on RBM to learn the relationship of face images between different modalities. Because RBM is nonlinear and unsupervised, our framework can learn the nonlinear relationship well and unlikely prone to overfitting.

The proposed framework includes 3 main steps: (1) extracting local Gabor features around facial points, as traditional face recognition methods do; (2) learning a shared representation by RBM for each group of local features; (3) processing the whole RBM representations by PCA and matching by Cosine similarity. Among them the key step is (2), in which a 3-layer RBM is constructed and the middle layer represents the shared properties of heterogeneous data.

The contributions of this paper are as follows.

1. A local to global learning framework is proposed for heterogeneous face recognition, which can achieve good results in all experiments.
2. For Sketch-Photo problem, perfect results are obtained on the CUFS [2] database. This is the first work that saturates the database.
3. Local RBMs are first used to learn the shared representations of heterogeneous face images. By plugging the local RBMs into the framework, we get state-of-the-art results on the CASIA HFB [10] and NIR-VIS 2.0 [11] databases.

2 Related Works

Heterogeneous face recognition research started from Tang and Wang's work in 2002 [12]. From that time to now, existing methods can be divided into two categories: Synthesis based and Classification oriented methods. In the early stage, the mainstream belongs to synthesis based methods, such as [13], [14] and [15]. [13] proposed a method, named as eigen-transformation, to synthesize photo by sketch and then recognized the identity in photo modality. To get more realistic results, [14] synthesized photo in a patch way, in which each image patch was first reconstructed by LLE and then stitched into a whole photo. [15]

also proposed a simple way to transform VIS to NIR face image. Although the results show that synthesis based method can achieve good visual quality, the recognition rate based on the synthesized images is moderate.

In late years, more classification oriented methods were proposed to improve the recognition rate directly. These methods just have one target: removing the difference of modalities, and meanwhile extracting discriminative feature. Many image processing and coding techniques are their essential parts, such as DoG filter [16], LBP, HOG [17], using which the difference between Sketch, NIR or VIS face images can be reduced significantly. Then, the processed heterogeneous data are mapped to a discriminative space by linear, nonlinear mapping [6][7] or random trees [18]. Because the target of this kind of methods is more direct than synthesis based methods, they always perform better.

Recently, several methods are proposed for multi-modal problems in deep learning community. [19] first proposed a multi-modal deep learning method based on denoising autoencoder, named as Bimodal Deep AE. But the Bimodal Deep AE performs poorly in Video-Audio matching experiments. On the contrary, another shallow architecture RBM-CCA results in surprisingly good performance. Unfortunately, [19] didn't give any analysis about why the deep net was worse than RBM-CCA. In 2012, [20] pointed out that in Bimodal Deep AE the responsibility of the multi-modal modeling fell entirely on the joint layer, and other layers gave no contributions. Therefore, they proposed a multi-modal Deep Boltzmann machine (DBM), which can spread out the responsibility of the multi-modal modeling over the entire network. Experiments illustrated the superiority of DBM in Image-Text retrieval task. Then, [21] applied the multi-modal DBM in the Image-Text retrieval challenge of ICML 2013 and got the first place in the challenge.

Because the multi-modal RBM in [20] has many good properties to deal with cross-modal matching problem, we plug the multi-modal RBM into the traditional face recognition pipeline to construct a novel framework for heterogeneous face recognition. By combing the advanced modules in these two fields, the proposed framework can work very well in challenging experiments.

3 Background

RBM has been widely used for modeling distribution of binary data. After Hinton's work [8], it became a standard building block of deep neural network. To model the real-valued data of face images, Gaussian RBM is used in this paper. This section will review the RBM, Gaussian RBM and Multi-modal RBM in brief.

3.1 Restricted Boltzmann Machines

RBM [9] is a generative stochastic graphical model that can learn the distribution of training data. The model consists of stochastic visible units $\mathbf{v} \in \{0, 1\}^m$ and

stochastic hidden units $\mathbf{h} \in \{0, 1\}^n$, which aims to minimize the following energy function:

$$E(\mathbf{v}, \mathbf{h}; \mathbf{a}, \mathbf{b}, \mathbf{W}) = -\mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}, \quad (1)$$

where \mathbf{a} is the biases of visible units; \mathbf{b} is the biases of hidden units; \mathbf{W} is the weights matrix to connect the visible and hidden units.

For image data, real-valued visible units $\mathbf{v} \in \mathbb{R}^m$ are used to replace the binary ones. The new model is called Gaussian RBM [22], the energy function of which is defined as:

$$E(\mathbf{v}, \mathbf{h}; \mathbf{a}, \mathbf{b}, \mathbf{W}) = \frac{1}{2} \mathbf{u}^T \mathbf{u} - \mathbf{b}^T \mathbf{h} - (\mathbf{v} \odot \frac{1}{\boldsymbol{\sigma}})^T \mathbf{W} \mathbf{h}, \quad (2)$$

where $\mathbf{u} = (\mathbf{v} - \mathbf{a}) \odot \frac{1}{\boldsymbol{\sigma}}$ denotes the normalized visible data. $\boldsymbol{\sigma}$ is a vector consisting of the standard deviations of each dimension. \odot denotes element-wise multiplication of vectors. Before training Gaussian RBM, the input data are usually normalized by WPCA or ZCA [23], *i.e.*, the standard deviations $\boldsymbol{\sigma}$ of the normalized data $\hat{\mathbf{v}}$ is 1. Then, the energy function can be simplified as:

$$E(\hat{\mathbf{v}}, \mathbf{h}; \mathbf{a}, \mathbf{b}, \mathbf{W}) = \frac{1}{2} (\hat{\mathbf{v}} - \mathbf{a})^T (\hat{\mathbf{v}} - \mathbf{a}) - \mathbf{b}^T \mathbf{h} - \hat{\mathbf{v}}^T \mathbf{W} \mathbf{h}. \quad (3)$$

Then the distribution over visible and hidden units is defined as:

$$P(\hat{\mathbf{v}}, \mathbf{h}; \boldsymbol{\theta}) = \frac{1}{Z} e^{-E(\hat{\mathbf{v}}, \mathbf{h}; \boldsymbol{\theta})}, \quad (4)$$

where $\boldsymbol{\theta}$ is an abberation for the parameters of RBM $\{\mathbf{a}, \mathbf{b}, \mathbf{W}\}$; Z is a partition function defined as the sum of $e^{-E(\hat{\mathbf{v}}, \mathbf{h}; \boldsymbol{\theta})}$ over all possible configurations.

3.2 Multi-modal RBM

[20] constructed a multi-modal RBM to model the relationship between image and text by combining a Gaussian RBM and Replicated Softmax RBM. For heterogenous face recognition problem, we use two Gaussian RBM to model the relationship between face data in two modalities. The structure of our model is shown in Figure 1. Its energy function is given by:

$$\begin{aligned} E(\hat{\mathbf{v}}_1, \hat{\mathbf{v}}_2, \mathbf{h}; \boldsymbol{\theta}) = & \frac{1}{2} (\hat{\mathbf{v}}_1 - \mathbf{a})^T (\hat{\mathbf{v}}_1 - \mathbf{a}) + \\ & \frac{1}{2} (\hat{\mathbf{v}}_2 - \mathbf{b})^T (\hat{\mathbf{v}}_2 - \mathbf{b}) - \\ & \mathbf{c}^T \mathbf{h} - \hat{\mathbf{v}}_1^T \mathbf{W}_1 \mathbf{h} - \hat{\mathbf{v}}_2^T \mathbf{W}_2 \mathbf{h}, \end{aligned} \quad (5)$$

where $\hat{\mathbf{v}}_1$ and $\hat{\mathbf{v}}_2$ are face images in two modalities; \mathbf{W}_1 and \mathbf{W}_2 are weights matrix for each modality respectively. The joint distribution over $\hat{\mathbf{v}}_1$, $\hat{\mathbf{v}}_2$, and \mathbf{h} can be calculated based on the energy function, as similar as Equ. (4).

Given the normalized training data $\hat{\mathbf{v}}_1$ and $\hat{\mathbf{v}}_2$, we can learn the parameters $\boldsymbol{\theta}$. Then, the trained multi-modal RBM can be used flexibly, such as

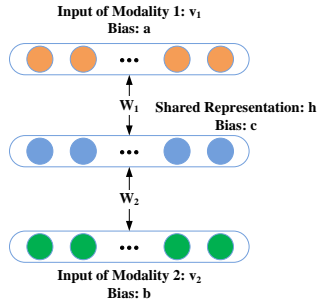


Fig. 1. A multi-modal RBM that modeling the joint distribution of face images in two modalities. The hidden layer in the model can be seen as a shared representation of the two input modalities.

1. generating missing modality by sampling from conditional distribution $P(\hat{\mathbf{v}}_1|\hat{\mathbf{v}}_2)$,
2. fusing two modalities by sampling from $P(\mathbf{h}|\hat{\mathbf{v}}_1, \hat{\mathbf{v}}_2)$,
3. inferring shared representation by sampling from $P(\mathbf{h}|\hat{\mathbf{v}}_1)$ and $P(\mathbf{h}|\hat{\mathbf{v}}_2)$ respectively.

Due to the experience in heterogeneous literature [6][7][18], this paper uses it for shared representation inference, which transforms the heterogenous data into a common space. For the details of multi-modal RBM learning and inference, please refer to [24][20].

4 Learning Shared Representation

4.1 Framework

The core of heterogeneous face recognition is modeling the relationship between different modalities and meanwhile reserving the discriminative information. To this end, we propose a framework for heterogeneous face recognition by incorporating RBM into the traditional face recognition pipeline. The flowchart of the framework is shown in Figure 2, in which the heterogeneous face images are illustrated by NIR and VIS for example. First, Gabor features are extracted at many facial points for two modalities respectively. Based on the Gabor features, a series of local RBMs are used to learn the shared representation of two modalities for each facial point. All local shared representations are then concatenated and processed by PCA. Finally the similarity of these modality-free features can be evaluated by Cosine metric.

The proposed framework has following advantages:

1. Local Gabor feature is the mainstream in face recognition, which has strong discriminative ability and is robust to variations;
2. We learn the shared representation locally because the modality gap is smaller in local region, and low dimensional data is more efficient for computation and easier to prevent overfitting;

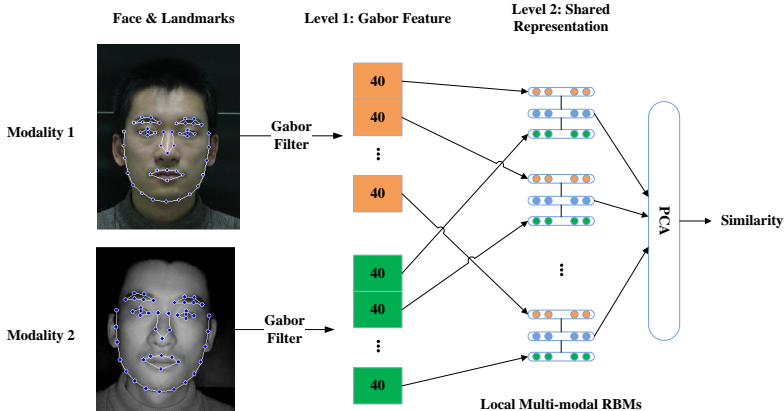


Fig. 2. The proposed framework for heterogeneous face recognition by combining traditional face recognition modules and local RBMs.

3. PCA can remove the redundance and heterogeneity further in holistic scale.

The details of each step in Figure 2 will be discussed in the following subsections.

4.2 Level 1 Representations

The task of level 1 is to extract discriminant and robust features for each modality. Recently, local features based on facial points achieved excellent performance in face recognition [25][26], especially in unconstrained face recognition, *e.g.*, LFW [27]. Although the face images in heterogeneous databases are both near frontal, facial points are still be used to deal with the small pose variations.

As shown in Figure 3, a standard set of facial points \mathbf{F}_s are defined for feature extraction and another 48 landmarks \mathbf{L}_s are defined for alignment, similar to [26]. Given a face image, we need put the facial points to the right place on it. [26] used a fast 3DMM model to do this work. For simplicity, this paper uses RBF warping [28] to transform the standard facial points to the face image. The warping process is shown in Figure 3. Given the landmarks \mathbf{L} of the input image, a warping function W can be solved based on \mathbf{L}_s and \mathbf{L} . Then the warped facial points are calculated by $\mathbf{F} = W(\mathbf{F}_s)$. We can see that the facial points can fit the input image well. The deformation factor of RBF warping is set to $0.1 \times$ “eye distance”.

At the warped 176×2 facial points \mathbf{F} , local features are extracted by a Gabor wavelet described in [29]. The space of Gabor wavelet is sampled in 8 orientations and 5 resolutions, thus giving $5 \times 8 = 40$ features for each facial point. Since the facial points are defined in a symmetric way, the features are grouped in left and right halves. Thus we get two feature vectors with 40×176 dimensions for each face image. Note that the facial symmetry trick has been used in many

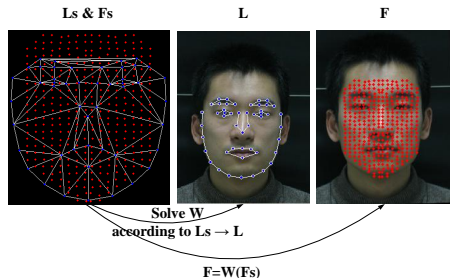


Fig. 3. The warping process of facial points. Left: Standard landmarks \mathbf{L}_s (blue dots) and facial points \mathbf{F}_s (red dots). Middle: A face image and its corresponding landmarks \mathbf{L} . Right: the warped facial points \mathbf{F} for the input image.

papers [30][10], which can augment the dataset and improve the computation efficiency.

4.3 Level 2 Representations

The task of level 2 is to build the relationship between two modalities. Previous work [31] has proven that the local relationship is easier to learn than holistic relationship, therefore we use local RBM to learn shared representation for each facial point. The structure of the RBMs is 40-80-40, including two input linear layers and a logistic hidden layer. Because the dimension of input of the RBM is very low, no sparse penalty and weight decay are used.

Existing methods, such as CSR [7], CITE [18] and their nonlinear versions, often learn the relationship in supervised and discriminative way. Different from them, RBM learns the joint distribution of the two modalities in a generative way, so RBM is less affected by overfitting. As described in [32], the relationships between modalities are not stationary with respect to the location in image, so we use many local RBMs to model them respectively, instead of one holistic RBM.

The level 1 features of two modalities are sent to 176 local RBMs, and their parameters are learned by using mean-field inference and an MCMC procedure described in [20]. In the training stage, the batch size is set to 10 and the number of batches is set to 50000. After the training is completed, we can infer the shared representations of two modalities by sampling from $P(\mathbf{h}|\hat{\mathbf{v}}_1)$ and $P(\mathbf{h}|\hat{\mathbf{v}}_2)$. While sampling from $P(\mathbf{h}|\hat{\mathbf{v}}_1)$, we treat $\hat{\mathbf{v}}_2$ and \mathbf{h} as missing data and initialize them randomly, then generate the hidden representation \mathbf{h} by alternating Gibbs sampler [20]. The hidden representation of another modality can be generated in a similar way. The activation probabilities of the hidden layer are called the shared representation of heterogeneous face images. The size of shared representation of a half face is 80×176 .

4.4 Cross-modal Matching

After the heterogeneity has been removed in local regions, the heterogeneity over holistic face still exists. As described in [11], PCA can capture the heterogeneity in its first several principle components, so we use PCA to process the feature in a holistic way. First, the 176 local representations are concatenated into a vector (the dimension is $80 \times 176 = 14080$) and the first several principle components of PCA are then removed. The number of removed components is tuned on the training set or development set. To this stage, the features of two modalities are actually transformed into a common space. Their similarity is calculated by Cosine metric. The similarity of two face halves are fused by sum rule.

We have also tried to learn a discriminative distance metric by LDA and Metric Learning based on the shared representations, but got worse results than PCA. The reason may be due to the limited data. We believe that the supervised methods will outperform PCA after having larger database in the future.

5 Experiments

To illustrate the superior performance of the proposed method, we take Sketch-Photo and NIR-VIS face recognition problems to conduct experiments. The results on three popular databases, CUFS [2], CASIA HFB [10] and CASIA NIR-VIS 2.0 [11], all outperform the current state-of-the-arts significantly.

5.1 Sketch to Photo

For Sketch-Photo problem, the CUFS database is used. The photos in CUFS come from three sources: 188 faces from CUHK student dataset, 123 faces from AR, and 295 faces from XM2VTS. Their corresponding sketches are drawn by an artist. In total, CUFS contains 606 subjects, 1 photo and 1 sketch per subject. As suggested in [18], the database is split into 306 training subjects and 300 testing subjects. To get unbiased results, the process repeats 10 times, and generates 10 splits. The mean and standard deviation of recognition rate of the 10 splits are reported. In the testing phase, photo is used as gallery and sketch is used as probe.

Because many good results have been reported on CUFS, it is considered as a relatively easy database [17]. For this simple experiment, the proposed method can work well without RBM, just using Gabor feature and PCA. Every photos and sketches are processed by facial points detection, Gabor feature extraction, PCA and Cosine matching. Only using Gabor feature, we can achieve comparable result to state-of-the-art methods, *i.e.*, Rank1=99.50%. By removing the first 20 principle components, the differences between photo and viewed sketch are removed completely, and enough identify information are reserved for classification. Without RBM, we get 100% recognition rate and outperform other compared methods. The comparisons are shown in Table 1.

Table 1. Rank1 recognition rates and VR@FAR=0.1% of various methods on CUFS.

	Rank1	VR
Gabor	99.50 \pm 0.39%	94.70 \pm 1.2%
Gabor + Remove 20 PCs	100 \pm 0%	100 \pm 0%
MRF+RS-LDA [2]	96.30%	N/A
LFDA [17]	99.47%	N/A
CITE [18]	99.87%	N/A

5.2 NIR to VIS

To illustrate the performance of our method further, two more difficult experiments are conducted on the CASIA HFB and NIR-VIS 2.0 databases.

CASIA HFB contains 2095 VIS and 3002 NIR face images from 202 subjects. We follow the evaluation protocol in [33] that selects 102 subjects for training and the other 100 subjects for testing. Similar to the previous experiment, the random selection is repeated 11 times. The first split (View 1) is used to tune the parameters of algorithm, and the other 10 splits (View 2) are used to report the performance.

CASIA NIR-VIS 2.0 is an upgraded version of HFB, the images in which are captured using the same devices as HFB, but has larger scale, contains more variations in pose, facial expression and age. CASIA NIR-VIS 2.0 is more close to practical applications. This database has standard evaluation protocols, so we use them directly.

In these two experiments, VIS face images are used as gallery and NIR face images are used as probe.

CASIA HFB First, the same framework in the previous experiment is applied on HFB, but just get moderate result. The reason may be that the first 20 principle components cannot capture the full difference between modalities. To verify this, we tune the number of removed principle components on View 1 carefully, but the VR (Verification Rate) cannot increase anymore. Thus we think the heterogeneity and discriminative information are coupled tightly and need to be dealt with in low level by RBM.

After introducing the RBM, the performance of our method increases significantly. As shown in Table 2, the VR@FAR=0.1% is improved from 71.70% to 92.25% and the deviation is also reduced remarkably. Meanwhile, the optimal number of removed principle components drops from 20 to 11 (see Figure 4), which indicates that the modality-free representations are successfully learned by local RBMs.

Compared to other methods in [33] and [34], the Rank1 and VR of our method are obviously higher. The SR (Sparse Representation) in [33] used the whole gallery to optimize the matching process, which has been proved can improve performance, especially in terms of ROC curve. For example, the VR of our method can be improved from 92.25% to 96.33% by using z-score normalization [35].

Because in face verification applications we cannot obtain the whole gallery, we just report the results without using the whole gallery. By fusing two classifiers and a commercial face recognition SDK, the VR of NN+SR+Cognitec [33] outperforms ours slightly. The reported performance of P-RS [36] is better than ours, but it is trained on larger training set. And P-RS is surely slower than our method because it's based on kernel similarities.

Table 2. Rank1 recognition rates and VR@FAR=0.1% of various methods on View 2 of CASIA HFB.

	Rank1	VR
Gabor	59.47±6.72%	33.51 ± 5.70%
Gabor + Remove 20 PCs	94.87 ± 1.72%	71.70 ± 6.42%
Gabor + RBM + Remove 11 PCs	99.38 ± 0.32%	92.25 ± 1.68%
NN [33]	88.8%	48.78 ± 3.87%
SR [33]	93.4%	77.56 ± 2.96%
NN+SR [33]	92.2%	79.05 ± 4.48%
Cognitec [33]	93.8%	85.62 ± 2.17%
NN+SR+Cognitec [33]	97.6%	93.45 ± 0.96%
C-DFD [34]	92.2%	65.5%
P-RS [36] ¹	-	95.8 ± 6.15%

Global, Convolutional and Local RBMs The layer in neural network has three popular styles: fully collected layer, locally collected layer with shared weights (convolutional) and locally collected layer with unshared weights (local). For RBM, we call them as global, convolutional and local RBMs. To illustrate the advantages of local RBMs, we plug them into our framework and compare their performances on View 1 of HFB, the information of which are shown in Table 3. The architecture of convolutional and local RBMs are all 40-80-40. Limited by the memory of our Geforce GTX670 GPU, the hidden layer of global RBM only uses 3520 units.

The complexity of the three kinds of RBMs are global > local > convolutional. Generally, complex models are easier to overfit to the training set and simple models are prone to underfitting. The results in Table 3 prove this point very well. The global RBM just performs well on the training set and the convolutional RBM performs moderately both on training and testing set. Among these models, the local RBMs obtain the best trade-off between complexity and generalization. Maybe the locality of connection and weight sharing can be finetuned further to get better results, but we leave this work to the future.

Parameter Tuning As discussed above, the number of removed principle components greatly affect the performance of our method. Generally, if the difference

¹ 133 subjects for training, 67 subjects for testing

Table 3. The comparison of Global, Convolutional and Local RBMs on View 1 of CASIA HFB. The 3rd column is VR@FAR=0.1% on the training set of View 1. The 4th column is VR@FAR=0.1% on the testing set of View 1.

	Architecture	VR (Train)	VR (Test)
Global	7040-3520-7040	99.94%	1.549%
Conv.	40-80-40	73.31%	71.79%
Local	$176 \times (40-80-40)$	99.45%	90.85%

between modalities is bigger, we need drop more principle components. However, there are also some identity information existing in these components, so we should find a trade-off. Figure 4 shows the relationship between the performance and the number of removed principle components on View 1. From the figure we can see that the performance of our method without RBMs are affected drastically by the first several principle components. But after using RBMs, the curves become smoother and quick to reach the optimal point. Therefore, we set the number of removed PCs to 20 when without RBMs and set the number to 11 when with RBMs.

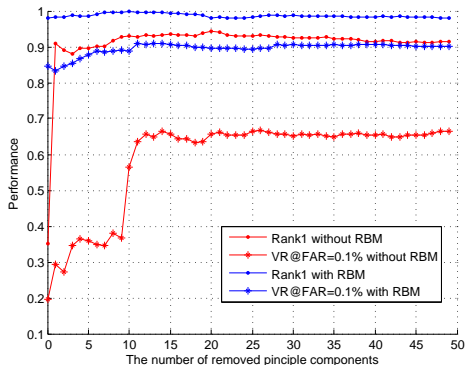


Fig. 4. The relationship between Rank1, VR and the number of removed principle components on View 1 of CASIA HFB. And the comparison curves of our method with/without RBMs.

Failure Cases Although the Rank1 recognition rate of our method is very high, there are still four failure cases on View 1 of HFB, which are shown in Figure 5. From the figure we can see that the four NIR probe images both have obvious variations in pose, specular reflectance on eyeglasses and expression. Even in traditional paradigm these factors heavily degrade the performance of face recognition, and they are more difficult to solve when coupling with spectrum variations.

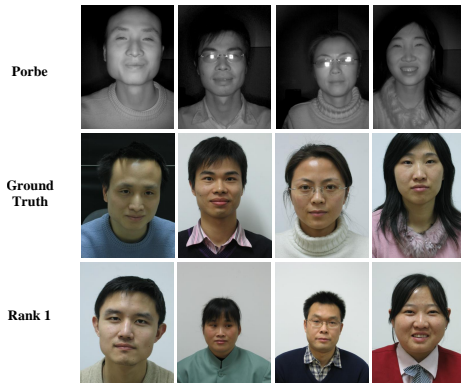


Fig. 5. The four failure cases on View 1 of CASIA HFB database. The first row are the NIR probe face images. The second row are the corresponding VIS face images of the first row. The third row are the retrieved Rank1 results of our method.

CASIA NIR-VIS 2.0 CASIA NIR-VIS 2.0 is a more challenging and practical database than the above two databases. The process of this experiment is as same as HFB, by first tuning the parameters on View 1 and then reporting results on View 2. From the results (Table 4) we can see that the Rank1 and VR on NIR-VIS 2.0 drop 10-20% compared to HFB. On this database, the improvements brought by removing the first PCs and RBM are still obvious, about 40% and 10% respectively. Because NIR-VIS 2.0 is a new database, we just list the baseline in [11] for comparison.

Table 4. Rank1 recognition rates and VR@FAR=0.1% of various methods on View 2 of CASIA NIR-VIS 2.0.

	Rank1	VR
Gabor	$36.18 \pm 2.56\%$	$33.37 \pm 2.29\%$
Gabor + Remove 20 PCs	$75.54 \pm 0.75\%$	$71.40 \pm 1.21\%$
Gabor + RBM + Remove 11 PCs	$86.16 \pm 0.98\%$	$81.29 \pm 1.82\%$
PCA+Sym+HCA [11]	$23.7 \pm 1.89\%$	19.27%

6 Conclusion

This paper proposed a novel framework for heterogeneous face recognition by combing RBM and the popular modules from face recognition. Because of its unsupervised nature, the framework is not prone to overfitting problem, and work well on many challenging heterogeneous face databases. Based on Gabor features, the modality-free shared representations were first learned successfully in low level by many local RBMs, and further processed by PCA in high level.

The proposed framework performed perfectly on the CUFS database and outperformed state-of-the-art methods significantly on the CASIA HFB and NIR-VIS 2.0 databases. Moreover, all experimental results illustrated the success of local RBMs to learn the shared representations. The future work will be conducted in two directions: (1) by stacking many multi-modal RBMs to learn high level representations; (2) exploring the way to fine tune the model with identity information.

References

1. Li, S.Z., Jain, A.K., eds.: The Encyclopedia of Biometrics. (2009)
2. Wang, X., Tang, X.: “Face photo-sketch synthesis and recognition”. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **31**(11) (2009) 1955–1967
3. Yi, D., Liu, R., Chu, R., Lei, Z., Li, S.Z.: “Face matching between near infrared and visible light images”. In: *Proceedings of IAPR International Conference on Biometric*, Seoul, Korea (August 2007)
4. Wang, X., Tang, X.: “Hallucinating face by eigentransformation”. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on* **35**(3) (2005) 425–434
5. Yan, J., Zhang, X., Lei, Z., Yi, D., Liao, S., Li, S.Z.: “Robust multi-resolution pedestrian detection in traffic scenes”. In: *IEEE Conference on Computer Vision and Pattern Recognition*. (2013)
6. Lin, D., Tang, X.: “Inter-modality face recognition”. In: *Proceedings of the European Conference on Computer Vision*. Volume 3954. (2006) 13–26
7. Lei, Z., Liao, S., Jain, A., Li, S.: “Coupled discriminant analysis for heterogeneous face recognition”. *IEEE Transactions on Information Forensics and Security* **7**(6) (2012) 1707–1716
8. Hinton, G.E., Salakhutdinov, R.R.: “Reducing the dimensionality of data with neural networks”. *Science* **313**(5786) (28 July 2006) 504–507
9. Smolensky, P.: *Parallel distributed processing: explorations in the microstructure of cognition*, vol. 1. MIT Press, Cambridge, MA, USA (1986) 194–281
10. Li, S., Lei, Z., Ao, M.: “The HFB face database for heterogeneous face biometrics research”. In: *IEEE Conference on Computer Vision and Pattern Recognition Workshops*. (2009) 1–8
11. Li, S.Z., Yi, D., Lei, Z., Liao, S.: “The CASIA NIR-VIS 2.0 face database”. In: *IEEE Conference on Computer Vision and Pattern Recognition Workshops*. (2013) 348–353
12. Tang, X., Wang, X.: “Face photo recognition using sketch”. In: *International Conference on Image Processing*. Volume 1. (2002) 257–260
13. Tang, X., Wang, X.: “Face sketch synthesis and recognition”. In: *IEEE International Conference on Computer Vision*. Volume 1. (2003) 687–694
14. Liu, Q., Tang, X., Jin, H., Lu, H., Ma, S.: “A nonlinear approach for face sketch synthesis and recognition”. In: *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Volume 1. (2005) 1005–1010
15. Wang, R., Yang, J., Yi, D., Li, S.: “An analysis-by-synthesis method for heterogeneous face biometrics”. In: *Proceedings of IAPR International Conference on Biometric*. (2009) 319–326

16. Liao, S., Yi, D., Lei, Z., Qin, R., Li, S.: “Heterogeneous face recognition from local structures of normalized appearance”. In Tistarelli, M., Nixon, M., eds.: *Advances in Biometrics*. Volume 5558 of *Lecture Notes in Computer Science*. (2009) 209–218
17. Klare, B., Li, Z., Jain, A.: “Matching forensic sketches to mug shot photos”. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **33**(3) (2011) 639–646
18. Zhang, W., Wang, X., Tang, X.: “Coupled information-theoretic encoding for face photo-sketch recognition”. In: *IEEE Conference on Computer Vision and Pattern Recognition*. (2011) 513–520
19. Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., Ng, A.Y.: “Multimodal deep learning”. In: *ICML*. (2011) 689–696
20. Srivastava, N., Salakhutdinov, R.: “Multimodal learning with deep boltzmann machines”. In: *Proceedings of Neural Information Processing Systems*. (2012) 2231–2239
21. Feng, F., Li, R., Wang, X.: “Constructing hierarchical image-tags bimodal representations for word tags alternative choice”. *CoRR* **abs/1307.1275** (2013)
22. Hinton, G.E.: “A practical guide to training restricted boltzmann machines”. In Montavon, G., Orr, G., Mller, K.R., eds.: *Neural Networks: Tricks of the Trade*. Volume 7700 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg (2012) 599–619
23. Bell, A., Sejnowski, T.J.: “The independent components of natural scenes are edge filters”. *Vision Research* **37** (1997) 3327–3338
24. Salakhutdinov, R., Hinton, G.E.: “Deep boltzmann machines”. *Journal of Machine Learning Research - Proceedings Track* **5** (2009) 448–455
25. Chen, D., Cao, X., Wen, F., Sun, J.: “Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification”. In: *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. (2013) 3025–3032
26. Yi, D., Lei, Z., Li, S.Z.: “Towards pose robust face recognition”. In: *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. (2013) 3539–3545
27. Huang, G.B., Ramesh, M., Berg, T., Learned-Miller, E.: “Labeled faces in the wild: A database for studying face recognition in unconstrained environments”. Technical Report 07-49, University of Massachusetts, Amherst (October 2007)
28. Arad, N., Reissfeld, D.: “Image warping using few anchor points and radial functions”. *Computer Graphics Forum* **14**(1) (1995) 35–46
29. Okada, K., Steffens, J., Maurer, T., Hong, H., Elagin, E., Neven, H., von der Malsburg, C.: “The Bochum/USC Face Recognition System and How it Fared in the FERET Phase III Test” (1998)
30. Berg, T., Belhumeur, P.: “Tom-vs-Pete classifiers and identity-preserving alignment for face verification”. In: *Proceedings of the British Machine Vision Conference*. (2012) 129.1–129.11
31. Yang, W., Yi, D., Lei, Z., Sang, J., Li, S.: “2D-3D face matching using cca”. In: *IEEE International Conference on Automatic Face Gesture Recognition*. (2008) 1–6
32. Huang, G., Lee, H., Learned-Miller, E.: “Learning hierarchical representations for face verification with convolutional deep belief networks”. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. (2012) 2518–2525
33. Klare, B., Jain, A.: “Heterogeneous face recognition: Matching NIR to visible light images”. In: *International Conference on Pattern Recognition*. (2010) 1513–1516
34. Lei, Z., Pietikainen, M., Li, S.: “Learning discriminant face descriptor”. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* (99) (2013)

35. Jain, A., Nandakumar, K., Ross, A.: “Score normalization in multimodal biometric systems”. *Pattern Recognition* **38** (2005) 2270–2285
36. Klare, B.F., Jain, A.K.: Heterogeneous face recognition using kernel prototype similarities. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **35**(6) (2013) 1410–1422