

CASCADE SUPPORT VECTOR REGRESSION-BASED FACIAL EXPRESSION-AWARE FACE FRONTALIZATION

Yiming Wang, Hui Yu, Honghai Liu

University of Portsmouth

Junyu Dong, Muwei Jian

Ocean University of China

ABSTRACT

The main aim of face frontalization is to synthesize the frontal facial appearances from non-frontal facial images. How to estimate the frontal face-shape is a crucial but very challenging problem in the frontalization task. Most existing methods use a single frontal face-template to fit in with frontal facial appearances, which will result in a loss of expression-related information. In this work, we present a novel facial expression-aware face frontalization method which directly learns the pair-wise relations between non-frontal face-shape and its frontal counterpart. The support vector regression is explored to train the pair-wise relation model. Considering non-linearity of the relationship, an appropriate cascade manner is applied to iteratively adjust and optimize the model. The frontal face-shape is then estimated via this model. With the estimated shape, frontal appearances are synthesized through a texture-fitting process formulated by solving a simple optimization problem. The proposed method has been evaluated on a in-the-wild facial expression database. The experimental results shows an outstanding performance of both facial expression-aware frontal face recovery and facial expression recognition.

Index Terms— Facial expression-aware face frontalization, cascade support vector regression, facial expression recognition

1. INTRODUCTION

Facial expression recognition (FER) in the wild addresses the unconstrained background and environment of facial images involving the challenging problem of large variations in head pose and occlusions [1]. FER based on 2D/3D view-invariant face-shape-free models have seldom been investigated. Most of view-invariant models focus on the problem of face recognition, but not FER. View-invariant face-shape-free FER requires accurate alignment of face shape between non-frontal face and its corresponding frontal face, which is always challenging under various facial expressions.

The existing view-invariant FER methods commonly focus on view modelling by learning the view-invariant features [2, 3] or training pose-wise classifiers [4, 5, 6, 7]. In order to ensure accuracy, most of the models need to be

trained per viewpoint/person/expression. Thus, a pose estimation step is always inevitable and the training data must be fairly large. Meanwhile, none of these methods address the problem of occlusions.

An appropriate way to overcome the above problems is to introduce face frontalization. Face frontalization is a comprehensive study which often involves face alignment, face morphing, face synthesizing etc. Substantial progress has been made on face recognition [8, 9, 10]. The main purpose of face frontalization is to recover the frontal facial appearances from unconstrained images. It often includes two steps: frontal face-shape estimation and frontal face-texture fitting.

Frontal face-shape estimation is the fundamental step of face frontalization. The existing work in [6, 7, 11, 12] either achieve only person-specific, not generic, frontalization or not able to recover the facial appearances. While frontal face-shape estimation is challenging, hard frontalization has shown its promising advantages. Currently, the approaches of [9] and [13] are the most effective generic face frontalization methods. These hard frontalization methods use a single 2D/3D reference facial template to fit facial appearances rather than estimating frontal face-shape. Both methods have successfully been tested for face recognition. However, hard frontalization may result in loss of expression-related cues since all the obtained frontal faces share a same template shape. Therefore, frontal face-shape estimation is inevitable for facial expression-aware face frontalization.

The task of texture-fitting is to compensate these missing parts of facial appearances. Current methods can be divided into interpolation-based approaches and symmetry-based compensation approaches. Most interpolation-based approaches are based on piece-wise affine warp that warps the non-frontal face-texture to frontal face-texture piece-wise [9, 14]. This strategy performs well on small head pose in terms of whatever pan or tilt angles. But it cannot deal with large head pose. On the contrary, symmetry-based compensation approaches can successfully recover the frontal face of large head pose in tilt angles, but will fail in pan angles [13].

In this paper, we focus on 2D generic frontalization and come up with an explicit idea of frontal-shape estimation and facial expression-aware face frontalization. Inspired by successful regression based 2D face alignment and face morphing methods in [10, 15, 16], we proposed a novel regression-

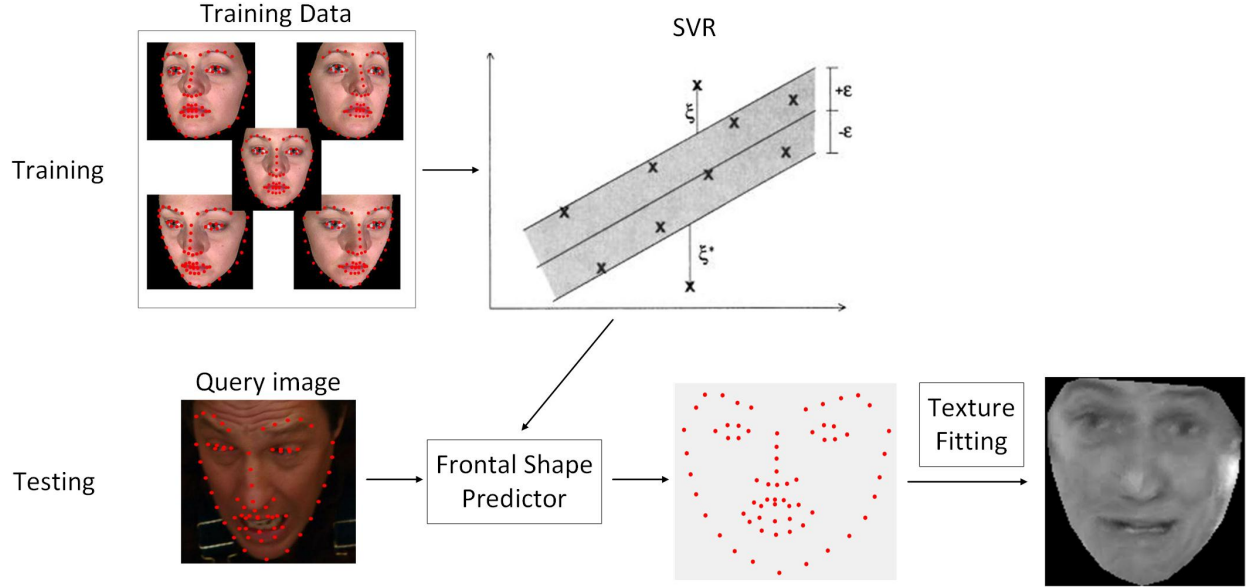


Fig. 1. Main process of proposed method

based method for 2D frontal face-shape estimation. As is shown in Fig 1, a pair of a non-frontal facial image and its corresponding frontal image are collected in order to train the regressor pair-wise. Support vector regression (SVR) is chosen due to its convexity and reasonability. Empirically, one step SVR cannot fit well with all shape instances because the changes in poses, expressions and individual characteristics are non-linear coupled. So we explored a cascade SVR to match with any possible shape features. In each cascade, we train a support vector regressor that best matches the groundtruth. During testing, the input face-shape will be through a sequentially cascade of SVR models.

After estimating the frontal face-shape, the input image will be warped from its original shape to the frontal one. We synthesize the frontal facial appearance by a linear combination of pre-defined eigen faces, where a group of parameters should be estimated. It, thus, become a optimization problem that minimize the differences between warped face and synthesized face. The derived frontal face will be in frontal view and robust to occlusions.

Without using reference template, the facial geometric model for each input image will be unique. The details of facial shape, individual features and facial expressions can be maintained after texture fitting. The main contribution include:

1) We propose a SVR-based framework for facial expression-aware face frontalization which is a new research branch of view-invariant FER. As far as we know, this is the first work in which template is totally withdrawn and an explicit idea of 2D frontal face shape estimation is proposed by using regression-based model.

2) Unlike many other view-invariant FER methods which

is extremely sensitive to the error of head pose estimation, there is no need for this method to estimate head pose, which effectively avoid the error accumulation.

3) Compared with traditional view-invariant FER, the proposed method is able to deal with occlusions.

4) The experimental results shows that this method outperforms the state-of-the-art FER approaches.

2. CASCADE SVR-BASED FACE FRONTALIZATION

2.1. Support Vector Regression

SVR is a supervised learning method for real number estimation. The main idea is to characterize the hyperplane that maximizes the margin. Given the training data $\{(x^1, y^1), \dots, (x^l, y^l)\}$, where $x \in \mathbb{R}^n, y \in \mathbb{R}$. Consider the linear function $f(x) = \langle \omega, x \rangle + b$, the SVR function can be expressed as:

$$\begin{aligned} \min \quad & \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^l (\xi_i^- + \xi_i^+) \\ \text{s.t.} \quad & \begin{cases} y_i - \langle \omega, x_i \rangle - b \leq \epsilon + \xi_i^+ \\ \langle \omega, x_i \rangle + b - y_i \leq \epsilon + \xi_i^- \\ \xi_i^+, \xi_i^- \leq 0 \end{cases} \end{aligned} \quad (1)$$

where $C > 0$ is a constant which make the balance between maximum margin and tolerance ϵ , and ξ is the slack variable which suggests that part of error is tolerated.

After introducing dual problem and Lagrangian multipli-

ers. The optimization problem becomes:

$$\begin{aligned} \max \quad & \begin{cases} \frac{1}{2} \sum_{i,j=1}^l (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle x_i, x_j \rangle \\ -\epsilon \sum_{i=1}^l (\alpha_i + \alpha_i^*) + \sum_{i=1}^l y_i (\alpha_i - \alpha_i^*) \end{cases} \quad (2) \\ \text{s.t.} \quad & \sum_{i,j=1}^l (\alpha_i - \alpha_i^*) = 0 \quad \text{and} \quad \alpha_i, \alpha_i^* \in [0, C] \end{aligned}$$

where α_i and α_i^* are the Lagrangian multipliers. The implementation of SVR varies. Quadratic programming is one of commonly used optimization method.

2.2. Frontal face-shape estimation

In this section, we present formulation the problem of SVR training. Given M annotated facial image pairs of non-frontal and corresponding frontal faces, the linear function of SVR can be defined as $\mathbf{x}_0 + \Delta \mathbf{x} = \langle \omega, \mathbf{x}_0 \rangle + b$, where \mathbf{x}_0 and \mathbf{x} represents the shape vector for non-frontal and frontal images, respectively. So $\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}_0$ is known. Consequently, the final function can be expressed as:

$$\Delta \mathbf{x} = \langle \omega, \mathbf{x}_0 \rangle + b \quad (3)$$

This equation can be referred to as the linear function of SVR.

Then we introduce the cascade manner of SVR regarding $\Delta \mathbf{x}^i$ representing the obtained $\Delta \mathbf{x}$ in the i th cascade. In each cascade, we revise Eq.3 into:

$$\Delta \mathbf{x}^i = \langle \omega^i, \mathbf{x}_0^i \rangle + b^i \quad (4)$$

and train a SVR model using:

$$\begin{aligned} \max \quad & \begin{cases} \frac{1}{2} \sum_{j,k=1}^l (\alpha_j^i - \alpha_j^{i,*})(\alpha_k^i - \alpha_k^{i,*}) k(x_{0,j}^i, x_{0,k}^i) \\ -\epsilon \sum_{j=1}^l (\alpha_j^i + \alpha_j^{i,*}) + \sum_{j=1}^l \Delta x_j^i (\alpha_j^i - \alpha_j^{i,*}) \end{cases} \quad (5) \\ \text{s.t.} \quad & \sum_{j,k=1}^l (\alpha_j^i - \alpha_j^{i,*}) = 0 \quad \text{and} \quad \alpha_j^i, \alpha_j^{i,*} \in [0, C] \end{aligned}$$

In the next cascade, we used the learned parameter to compute

$$\Delta \mathbf{x}^i = \sum_{j,k}^l (\alpha_j^i - \alpha_j^{i,*}) k(x_j^i, x_k^i) \quad (6)$$

In the new round, $\mathbf{x}_0^{i+1} = \Delta \mathbf{x}^i + \mathbf{x}_0^i$, $\Delta \mathbf{x}^{i+1} = \mathbf{x} - \mathbf{x}_0^{i+1}$ and they will be used to train this round SVR model.

With the cascade SVR iteratively moving on, the value of the parameters, C and ϵ , will gradually reduce. The algorithm will stop until these parameters and $\Delta \mathbf{x}^i$ turn zero. The cascade SVR empirically needs 4 to 5 steps.

During testing, the non-frontal facial landmarks should be localized first. There are many existing facial landmark detection methods that has been proved to be effective. With the obtained facial landmarks, frontal face-shape will be estimated using Eq. 6 sequentially.

2.3. Face-texture fitting

Let $I \in \mathbb{R}^{m \times n}$ describe a non-frontal facial image whose landmarks is \mathbf{x}_0 and corresponding frontal shape is \mathbf{x} . The piece wise affine warp strategy can be directly used to compute the warped image $\mathbf{W}(I; p)$ [17].

The frontal facial appearances are synthesized via a linear combination of a set of pre-defined eigen faces:

$$F = Uc \quad (7)$$

where U is eigen face vector and c is parameters. Consequently, the optimization problem can be solved by:

$$\arg \max_c \|W(I; p) - Uc\|^2 \quad (8)$$

With a simple derivation, c can be obtained by $c = UW(I; p)$. The synthesized frontal face can be obtained by Eq. 7.

3. EXPERIMENT

The performance of the proposed method has been validated in two tasks: 1) frontal face reconstruction 2) FER in the wild.

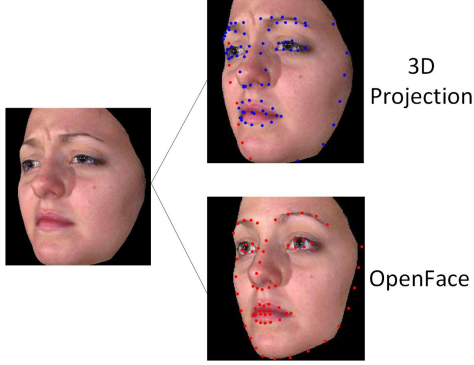
3.1. Prepare for training data

Binghamton University 3D Facial Expression (BU3DFE) is a static 3D facial expression database which include 100 subjects with 2500 3D facial expression models. The training data are captured by rendering 2D images using 3D models. Images are captured at 7 pan angles (-45°, -30°, -15°, 0°, 15°, 30°, 45°) and 5 tilt angles (-30°, -15°, 0°, 15°, 30°), which results in totally 35 different viewpoints. Each training instance includes the position landmark points in one of the 34 non-frontal rotations and the corresponding points in frontal pose.

BU3DFE provide the 3D position of 83 landmarks for each 3D facial model. When the 3D landmark points were projected to 2D plane, there would be misalignments especially when there was large out-of-plane rotation. As is shown in figure 2, the red points of 3D projection are obviously misaligned. So we use OpenFace to automatically detect landmark points. OpenFace [19] is a very simple and effective tool for facial landmark detection. Most landmarks can be well detected using this software. Misaligned points were manually revised.

Table 1. Recognition rate (%) of different methods on SFEW database

	Angry	Disgust	Fear	Happy	Neutral	Sadness	Surprise	Total
Baseline	23.00	13.00	13.90	29.00	23.00	17.00	13.50	18.90
[6]	25.89	28.24	17.17	42.98	14.00	33.33	10.99	24.70
[18]	24.11	14.12	20.20	50.00	23.00	23.23	21.98	26.14
Proposed	40.18	25.88	48.48	55.26	37.00	36.36	37.36	40.71

**Fig. 2.** Two different ways of collecting training data

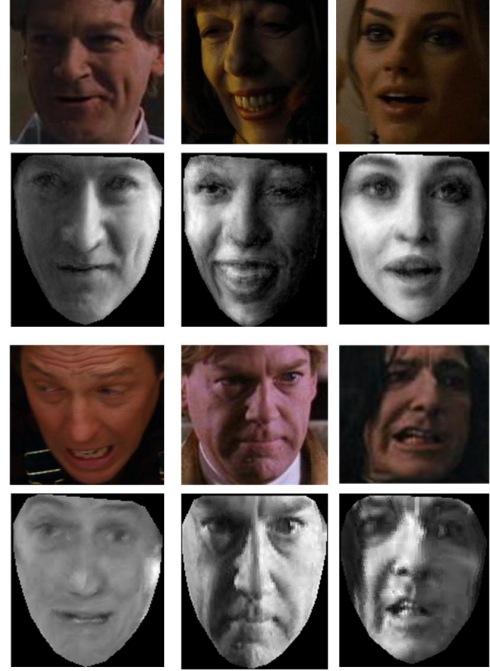
3.2. Frontalization and recognition

The shape model can be trained on the rendered 2D images from BU3DFE. For the two parameters of SVR, C is set by 10, 1, 0.1, 0.01 at the four iterations and ϵ is fixed to 0.

In order to evaluate the performance on the unconstrained images, we use another database for testing. Statistical Facial Expression in the Wild (SFEW) [20] is a spontaneous facial expression database. It contains 700 images captured from movies labelled by seven categories: six universal emotions and neutral.

From Figure 3 we can see the outstanding performance of face frontalization. Whatever pan rotation or tilt rotation can be well recovered to the frontal view. Meanwhile, the facial expression related cues are maintained.

For FER, there is a standard evaluation protocol provided by the authors of SFEW. The evaluation is strictly person-independent. In this experiment, Local Binary Pattern (LBP) and Support Vector Machine (SVM) are used for feature extraction and emotion classification, respectively. In table 1, the methods of [6] and [18] are small sample learning methods which currently achieved the best results. It is obvious that our method outperforms the state-of-the-art approaches. The overall recognition rate of the proposed method is 10% higher than the others, which suggests a considerable improvement. Based on this result we can conclude that facial expression-aware face frontalization can achieve promising result for FER in the wild. In this comparison, we did not

**Fig. 3.** Examples of face frontalization

mention deep learn because our work focus on small sample learning which is quite different from deep learning. Meanwhile, deep learning methods must use a large volume of external training data, which is not totally agree with the evaluation protocol of SFEW.

4. CONCLUSIONS

In this paper, we have presented a novel regression-based frontal face-shape estimation method for facial expression-aware face frontalization and applied it to FER in the wild. As far as we know, this is the first work that can achieves both face frontalization and facial expression recovery. The experiment on unconstrained images shows impressive visual effects of the synthesized faces and a significant improvement in facial expression recognition for the data in the wild.

5. REFERENCES

- [1] M. Valstar, B. Jiang, M. Mehu, M. Pantic, and K. Scherer, "The first facial expression recognition and analysis challenge," in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*, pp. 921–926, 2011.
- [2] U. Tariq, J. Yang, and T. S. Huang, "Multi-view facial expression recognition analysis with generic sparse coding feature," in *Proceedings of European Conference on Computer Vision*, pp. 578–588, 2012.
- [3] U. Tariq, J. Yang, and T. S. Huang, "Supervised super-vector encoding for facial expression recognition," *Pattern Recognition Letters*, vol. 46, pp. 89–95, 2014.
- [4] S. Moore and R. Bowden, "Local binary patterns for multi-view facial expression recognition," *Computer Vision and Image Understanding*, vol. 115, no. 4, pp. 541–558, 2011.
- [5] N. Hesse, T. Gehrig, H. Gao, and H. K. Ekenel, "Multi-view facial expression recognition using local appearance features," *International Journal of Computer Vision*, vol. 83, no. 2, pp. 178–194, 2011.
- [6] S. Eleftheriadis, O. Rudovic, and M. Pantic, "Discriminative shared gaussian processes for multiview and view-invariant facial expression recognition," *IEEE Transaction on Image Processing*, vol. 24, no. 1, pp. 189–204, 2015.
- [7] O. Rudovic, M. Pantic, and I. Patras, "Coupled gaussian processes for pose-invariant facial expression recognition," *IEEE Transaction on Pattern Analysis and Machine Intelligence*, vol. 35, no. 6, pp. 1357–1369, 2013.
- [8] G. Tzimiropoulos and M. Pantic, "Fast algorithms for fitting active appearance models to unconstrained images," *International Journal of Computer Vision*, pp. 1–17, 2016.
- [9] C. Sagonas, Y. Panagakis, S. Zafeiriou, and M. Pantic, "Robust statistical face frontalization," in *Proceedings of IEEE International Conference on Computer Vision*, pp. 3871–3879, 2015.
- [10] A. Akshay, S. Lucey, and R. Goecke, "Regression based automatic face annotation for deformable model building," *Pattern Recognition*, vol. 44, no. 10, pp. 2598–2613, 2011.
- [11] L. A. Jeni, J. F. Cohn, and T. Kanade, "Dense 3d face alignment from 2d videos in real-time," in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*, pp. 1–8, 2015.
- [12] J. Roth, Y. Tong, and X. Liu, "Unconstrained 3d face reconstruction," in *IEEE International Conference Computer Vision Workshops*, pp. 2606–2615, 2015.
- [13] T. Hassner, S. Harel, E. Paz, and R. Enbar, "Effective face frontalization in unconstrained images," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4295–4304, 2015.
- [14] H. T. Ho and R. Chellappa, "Pose-invariant face recognition using markov random fields," *IEEE Transaction on Image Processing*, vol. 22, no. 4, pp. 1573–1584, 2013.
- [15] X. Xiong and F. D. la Torre, "Supervised descent method and its applications to face alignment," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 532–539, 2013.
- [16] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. Zafeiriou, "Mnemonic descent method: A recurrent process applied for end-to-end face alignment," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [17] I. Matthews and S. Baker, "Active appearance model revisited," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 135–164, 2004.
- [18] M. Liu, S. Li, and X. Chen, "Au-aware deep networks for facial expression recognition," in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*, pp. 1–6, 2013.
- [19] T. Baltruaitis, P. Robinson, and L. Morency, "Openface: an open source facial behavior analysis toolkit," in *IEEE Winter Conference on Applications of Computer Vision*, pp. 1–10, 2016.
- [20] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon, "Static facial expression analysis in tough conditions: Data evaluation protocol and benchmark," in *IEEE International Conference Computer Vision Workshops*, pp. 2106–2112, 2011.