

A TWO-STEP LEARNING METHOD FOR DETECTING LANDMARKS ON FACES FROM DIFFERENT DOMAINS

Bruna Vieira Frade

Erickson R. Nascimento

Universidade Federal de Minas Gerais (UFMG), Brazil
{brunafrade, erickson}@dcc.ufmg.br

ABSTRACT

The detection of fiducial points on faces has significantly been favored by the rapid progress in the field of machine learning, in particular in the convolution networks. However, the accuracy of most of the detectors strongly depends on an enormous amount of annotated data. In this work, we present a domain adaptation approach based on a two-step learning to detect fiducial points on human and animal faces. We evaluate our method on three different datasets composed of different animal faces (cats, dogs, and horses). The experiments show that our method performs better than state of the art and can use few annotated data to leverage the detection of landmarks reducing the demand for large volume of annotated data.

Index Terms— Landmarks detection, Machine Learning, Domain Adaptation, Human faces, Animal faces

1. INTRODUCTION

Detecting landmarks embedded with semantic information from images is one of the key challenges in image processing and computer vision fields. In general, landmarks or fiducial points are related to discriminative locations in the image, frequently embedding some meaning. For example, in human or animal faces, a landmark locates regions comprising the eyes, eyebrows, mouth, and the tip of the nose. After all, the automatic estimation of landmarks on faces has a myriad of applications such as faces recognition, game animation, avatars, and transferring facial expressions [1].

Despite remarkable advances in detecting landmarks on human faces, most of the methods require a large number of annotated data. For every different type of face like an animal such as a cat or a dog face, we still have to use the time-consuming process of annotating each landmark for a considerable amount of data. In other words, although detecting the same type of landmarks present in a large dataset of human faces (*e.g.*, eyes, nose, *etc.*), we need to build an entirely new dataset. Thus, a central challenge in facial landmark detection is how to use the annotation available in big datasets such those for human faces to improve the detection of similar landmarks but on different types of faces.

In this work, we present a domain adaptation algorithm based on deep learning to detect landmarks on human and non-humans faces. Our method builds a landmark detector by performing two tasks: i) learning to identify landmarks in a supervised way by using labeled data of human faces (source domain) and ii) learning to reconstruct non-human faces with unlabeled data (target domain). The final representation of our method preserves the discriminability from the labeled data and encodes the landmarks locations of the target domain. This capability suggests that the performance of our method stems from the creation of a single representation that encodes the structure information of non-human face and relevant features for the landmark detection on the human face.

According to our experiments, our method outperformed the state-of-the-art landmarks detector for interspecies face [2] up to 10% precision gap. We evaluated our method on a variety of type of faces, where it was capable of learning the cross-domain landmark detection task, without requiring a big collection of annotated data.

Related Work. In the past several years, a popular approach for detecting fiducial points was based on using classifiers [3, 4, 5]. However, recently we have witnessed an explosion of approaches to learning features found on convolution neural networks. Sun et al. [6] used three levels of convolutional neural networks (CNN) to estimate the position of landmarks on faces. Thanks to the high-level global characteristics extracted from the entire face, their method increases the precision of the landmarks detection. To minimize occlusion effects, Zhang et al. [7] proposed a multitask learning to optimize landmark detection through heterogeneous but correlated tasks, *i.e.*, head pose and facial attributes inference. Zhang et al. [8] aim at refining the alignment in each stage. Using an autoencoder approach, they tried to predict landmarks quickly through low-quality images and progressively improving previous results while increasing its resolution. The work of Yu et al. [9] used a cascading approach called deep deformation network (DDN). The DDN learns to extract the shape information and then uses a landmark transformation network to estimate the local parameterized defor-

mation aiming to refine first step results.

A significant disadvantage of most CNN based methods is the need for large-scale datasets. Yang et al. [10] and Rashid et al. [2] propose to adjust the landmark learning from human faces to other target topology. The work of Rashid et al. searches for similar human faces for each animal sample in an unsupervised method. Yang et al. interpolate face characteristics by cascade regression aiming at detecting through shape using fewer samples data of the target domain. A recent approach is the Deep Reconstruction Classification Network (DRCN) [11]. The DRCN jointly learns a shared encoding representation for the digits classification task. Inspired by the DRCN approach, our method is also based on a two-step learning domain adaptation approach, but different from it, we used the learned encoding for a regression problem solution in the faces domain.

2. METHODOLOGY

Let \mathcal{D}^{source} be a large set of labeled images like human faces dataset, and \mathcal{D}^{target} be a set with a small number of labeled samples, *e.g.*, dog’s face. Our methodology has been designed to detect landmarks in both \mathcal{D}^{source} and \mathcal{D}^{target} domains. Our formulation is based on a two-step learning approach, wherein the first step it learns to reconstruct images from \mathcal{D}^{target} using an unsupervised strategy and in the second step it solves a regression problem in a supervised way predicting the coordinates of the landmarks in faces from \mathcal{D}^{source} . Figure 1 illustrates the whole process.

The reconstruction and detection steps rely on the network $g_{enc}(\cdot)$ that learns to encode discriminative features of faces on different domains and tasks. This function plays a key role in building a model capable of reconstructing faces and detecting landmarks.

Reconstruction. Let $\mathbf{u}_i \in \mathcal{D}^{target}$ be the i -th unlabeled image of an animal. After encoding the face features using $g_{enc}(\mathbf{u}_i)$, we apply a decode function $g_{dec}(\cdot)$ in order to create an output $g_{dec}(g_{enc}(\mathbf{u}_i))$ as close as possible to \mathbf{u}_i .

In other words, we train the network to minimize the loss function:

$$\mathcal{L}_{rec} = \frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_i - g_{dec}(g_{enc}(\mathbf{u}_i))\|^2, \quad (1)$$

where n is the number of images in \mathcal{D}^{target} with no annotations, because this step does not require labeled data.

Landmark detection. In the second step, we apply a supervised approach to learn to detect landmarks on faces using annotated data. Let $\mathbf{x}_i \in \mathcal{D}^{source}$ be the i -th image of a person and $\mathbf{y}_i \in \mathbb{R}^6$ be the image coordinates of the landmarks, *i.e.*, eyes and nose.

Algorithm 1 Two-step learning for landmark detection.

```

1: Labeled dataset:  $\mathcal{D}^{source} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$ 
2: Unlabeled dataset:  $\mathcal{D}^{target} = \{\mathbf{u}_j\}_{j=1}^n$ 
3: for each  $e < totalEpoch$  do
4:   for each  $batch_t \in \mathcal{D}^{target}$  and  $batch_s \in \mathcal{D}^{source}$  do
5:     ForwardUnsupervised( $batch_t$ )
6:     ComputeReconstructionError:  $\mathcal{L}_{rec}(batch_t)$ 
7:     UpdateWeights of  $g_{dec}$  and  $g_{enc}$ 
8:     ForwardSupervised( $batch_s$ )
9:     ComputeRegressionError:  $\mathcal{L}_{reg}(batch_s)$ 
10:    UpdateWeights of  $g_{reg}$  and  $g_{enc}$ 
11:   end for
12: end for

```

In our approach, the supervised branch of our architecture uses the representation learned in the unsupervised step for feeding the supervised the regression function $g_{reg}(\cdot)$. The regression function learns the coordinates of each landmark by computing the error between the error between the ground truth \mathbf{y}_i and the prediction $g_{reg}(g_{enc}(\mathbf{x}_i))$ as:

$$\mathcal{L}_{reg} = \frac{1}{m} \sum_{i=1}^m \text{MAE}(\mathbf{y}_i - g_{reg}(g_{enc}(\mathbf{x}_i))), \quad (2)$$

where m is the size of the annotated set and MAE is the mean absolute error between two vectors of size k , *i.e.*, $\text{MAE}(\mathbf{a}, \mathbf{b}) = \frac{1}{k} \sum_{i=1}^k |\mathbf{a}_i - \mathbf{b}_i|$.

In our two-step learning, we address the lack of annotated data in the target domain. Algorithm 1 depicts the learning loop. Instead of starting the learning from random weights, we use the patterns learned from a reconstruction task. These patterns help to leverage the landmark detection when solving the regression in a supervised way but in a different domain. Thanks to this strategy, we can extract robust features by hallucinating face features of datasets with a few annotations. Our hypothesis is grounded in the idea that we can learn a function that maps similar features in both domains. Thus, with this mapping, we can simultaneously perform the regression and reconstruction tasks and simultaneously learn how to detect landmarks in the target domain.

Unlike DRCN [11] that updates the weights of a classification net using all labeled batches and then adjusts the encoder weights in the reconstruction with all unlabeled batches, our approach updates the weights of the reconstruction and regression nets for each batch intercalating both steps. This fact plays a key role in the development of a robust and flexible strategy that can work with unlabeled data in the target domain or a small collection of annotations.

3. EXPERIMENTS

Datasets. In our experiments, the labeled data are from human faces of the Keggler [12]. For the target domain, *i.e.*, an

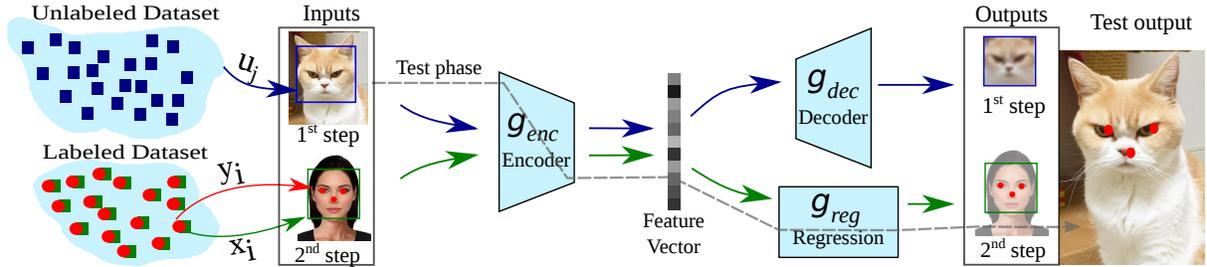


Fig. 1. Illustration of our two-step learning with the supervised (regression) and the unsupervised steps. In the first step, we encode the face features of a unlabeled image u_i using g_{enc} network. Then we apply a decode function to reconstruct the input, *i.e.*, $g_{dec}(g_{enc}(u_i))$. In the second step, we feed a regression function g_{reg} that learns the coordinates of each landmark by estimating the error between the ground truth y_i and the prediction $g_{reg}(g_{enc}(x_i))$.

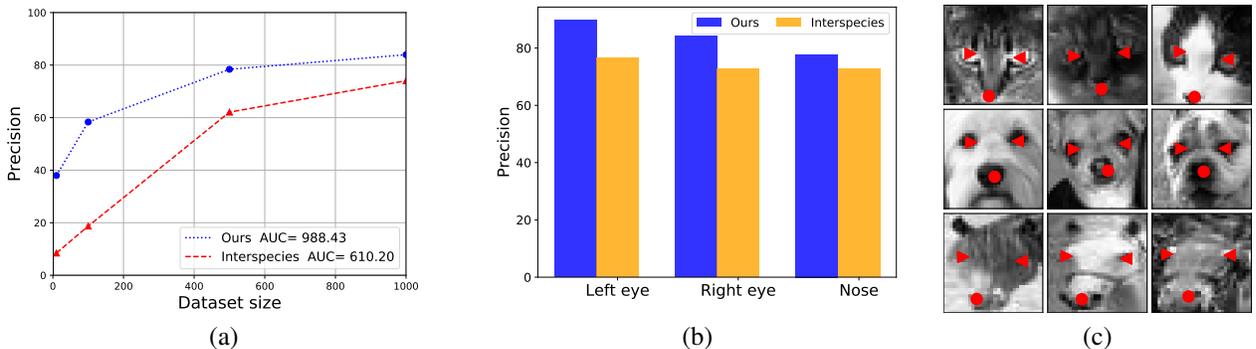


Fig. 2. (a) ROC curve of our method and the Interspecies using different numbers of labeled images of the target domain (from 10 to 1,000 images). Larger AUC means better performance. (b) Precision of the predicted landmarks using with 1,000 of labeled images. (c) Qualitative results for cats, dogs, and horses training our method with 100 of labeled images.

imal faces, aside from dataset of cat faces [13], we also evaluated our approach on dogs faces [14] and horse faces [2]. All datasets have the faces annotated with 3 landmarks: left and right eyes, and nose. We used the Euclidean distance as evaluation metric for the location predictions and margin of error of 10% of the image size (in the dataset we used in the experiments corresponds to a radius of 3 pixels) to classify a detection as correct. We performed data augmentation in source and target datasets by applying a random rotations in the images from -30 to 30 degrees. Moreover, we applied translations and noise in the target set to ensure robust results in the reconstruction step.

Baselines. We compared our method with a standard convolutional network (ConvNet) for supervised landmark detection. This ConvNet has the same architecture of our supervised net and it was trained on the target domain with labeled data. We also pit our detector against the state-of-the-art method on landmark detection for faces in different topologies called Interspecies [2]. We used the code provided by the authors with our configuration of training and test data ranging from 10 to 1,000 images.

Implementation. The encoder $g_{enc}(\cdot)$ has five convolution layers with 3×3 filters and padding 1×1 : 300 filters in conv1, 250 filters in conv2, 200 filters in conv3, 150 filters in conv4, and conv5 with 100 filters; two pooling layers 2×2 after the first and second convolutional layers (pool1 and pool2) and a fully connected layer (fc4) with 500 neurons. The decoder net g_{dec} is the mirror image of the encoder architecture. In the regression, we feed a fully connected layer (fc-regressor) with the output from the g_{enc} net. We used ReLU in all hidden and output layers and hyperbolic tangent activation for the regression layer. We used a learning rate equals to 3×10^{-4} for point detection and the reconstruction. In the training, we ran 500 epochs and used a batch size equals to 128. The size of all input images is 32×32 . All source code and experimental data will be publicly available.

3.1. Comparison against the state of the art

We compared our work with the Interspecies method proposed by Rashid et al. [2], the current state of the art in detecting landmarks on animal faces. In the experiments, we varied

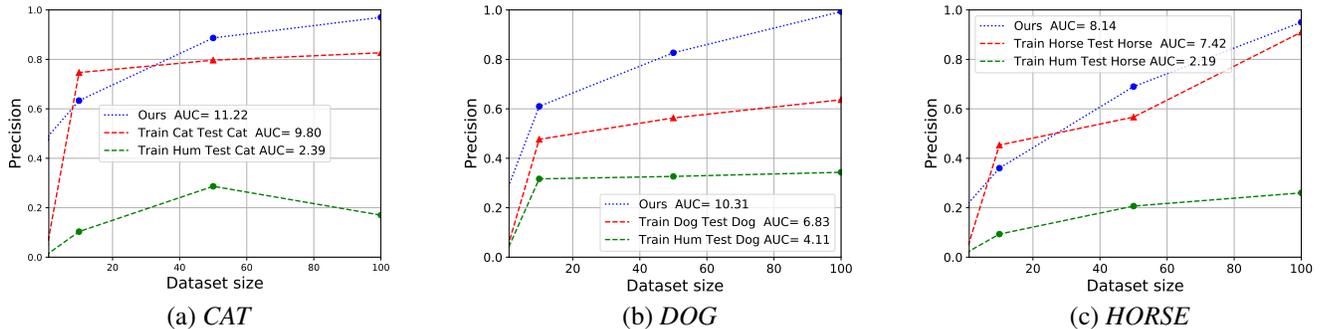


Fig. 3. ROC curve varying the number of labeled data for CAT, DOG and HORSE datasets. Our method (blue curve) was superior than ConvNet trained with samples from target domain (red curve) and source domain (green curve).

Table 1. Area Under the Curve (AUC) for ROC of precision for our method and a ConvNet trained with labeled data from a human faces dataset. Best in bold.

Method	Dataset (AUC)		
	CAT	DOG	HORSE
OURS	86.3	82.3	79.97
CONVNET	77.33	81.10	76.65

the size of training from 10 to 1,000 images.

Figure 2 (a) shows the Receiver Operating Characteristic (ROC) of our method and the Interspecies. One can clearly see that our method outperformed the Interspecies method. Our method obtained 83% of the precision while the Interspecies method obtained 73.2% of correct detections. We also evaluate prediction precision for each landmark. As one can see in the bars of Figure 2 (b), our method has larger precision for all the landmarks, *i.e.*, nose and eyes locations.

From these results, we can draw the following observations. First, even when there is no labeled data from the target domain, our method performs better than Interspecies. The experiments show that our approach can hallucinate features by adapting learnt features in a reconstruction task from a different domain. Second, whenever available, our method can use few annotated data to leverage the detection of landmarks reducing the demand of large volume of annotated data.

3.2. Ablation analysis

For a more detailed performance assessment, we also evaluate our method in two experiments: in the first experiment we used only unlabeled data from the target domain and labeled data from the source domain; in the second experiment, we gradually increased the size of labeled dataset from the target domain starting from 5 up to 100 images.

Only unlabeled data from the target domain. Table 1

shows the area under the curve of the ROC curves for our model trained with no labeled data from target domain and a ConvNet trained with source domain only. When comparing with ConvNet, our method improved the detection in all three datasets. Together these results show that the hypothesis of leveraging feature from one domain to another holds for detecting landmarks on faces.

Using a few of labeled data from the target domain. In this experiment, we use some labeled data from the target domain in the regression step. The idea is to analyze the performance of our approach when using a small number of labeled data from the target domain. For each batch, we forward a subset of labeled data from the target domain and solve the regression according to the error between the prediction and the ground truth of the target domain. Then, the network weights are updated. We used four different number of labeled images: 0, 10, 50, and 100. Figure 3 shows the precision for the three datasets. It is noteworthy the rapid increase in the precision of our method and the superior performance (larger AUC). This fact can be explained by the transfer learning from the source domain and the features learning in the facial reconstruction in the target domain.

4. CONCLUSION

We presented a novel method for detecting landmarks on faces in different domains such as human and animal faces. Our method is based on a two-step learning (supervised and unsupervised) that reduces the need for a large annotated data. The experiments show that our method performed better than the state-of-the-art method even when there is a small collection of annotated data from the target domain.

5. ACKNOWLEDGMENTS

The authors would like to thank the agencies CAPES, CNPq, FAPEMIG, Vale Institute of Technology (ITV), and Petrobras for funding different parts of this work.

6. REFERENCES

- [1] Robert W Sumner and Jovan Popović, “Deformation transfer for triangle meshes,” *ACM Transactions on Graphics (TOG)*, vol. 23, no. 3, pp. 399–405, 2004. [1](#)
- [2] Maheen Rashid, Xiuye Gu, and Yong Jae Lee, “Interspecies knowledge transfer for facial keypoint detection,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [1](#), [2](#), [3](#)
- [3] Lin Liang, Rong Xiao, Fang Wen, and Jian Sun, “Face alignment via component-based discriminative search,” *European Conference o Computer Vision (ECCV)*, pp. 72–85, 2008. [1](#)
- [4] Xiangxin Zhu and Deva Ramanan, “Face detection, pose estimation, and landmark localization in the wild,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 2879–2886. [1](#)
- [5] Peter N Belhumeur, David W Jacobs, David J Kriegman, and Neeraj Kumar, “Localizing parts of faces using a consensus of exemplars,” *IEEE transactions on pattern analysis and machine intelligence (TPAMI)*, vol. 35, no. 12, pp. 2930–2940, 2013. [1](#)
- [6] Yi Sun, Xiaogang Wang, and Xiaoou Tang, “Deep convolutional network cascade for facial point detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 3476–3483. [1](#)
- [7] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang, “Facial landmark detection by deep multi-task learning,” in *European Conference o Computer Vision (ECCV)*. Springer, 2014, pp. 94–108. [1](#)
- [8] Jie Zhang, Shiguang Shan, Meina Kan, and Xilin Chen, “Coarse-to-fine auto-encoder networks (cfan) for real-time face alignment,” in *European Conference o Computer Vision (ECCV)*. Springer, 2014, pp. 1–16. [1](#)
- [9] Xiang Yu, Feng Zhou, and Manmohan Chandraker, “Deep deformation network for object landmark localization,” in *European Conference o Computer Vision (ECCV)*. Springer, 2016, pp. 52–70. [1](#)
- [10] Heng Yang, Renqiao Zhang, and Peter Robinson, “Human and sheep facial landmarks localisation by triplet interpolated features,” in *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*. IEEE, 2016, pp. 1–8. [2](#)
- [11] Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, David Balduzzi, and Wen Li, “Deep reconstruction-classification networks for unsupervised domain adaptation,” in *European Conference o Computer Vision (ECCV)*. Springer, 2016, pp. 597–613. [2](#)
- [12] “Facial keypoint detection competition,” Kaggle. [2](#)
- [13] Jian Sun Weiwei Zhang and Xiaoou Tang, “Cat head detection - how to effectively exploit shape and texture features,” *European Conference o Computer Vision (ECCV)*, 2005. [3](#)
- [14] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei, “Novel dataset for fine-grained image categorization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Colorado Springs, CO, June 2011. [3](#)