

GANs for Generating Different Face Expressions

Helena Čeović

Josipa Lipovac

Antun Magdić

Andrea Omićević

Ante Žužul

Abstract—Write abstract here.

Index Terms—StrGAN...

I. INTRODUCTION

There have been several researches that tackled the problem of facial expression recognition. However, even though the problem of facial expression generation is more challenging than expected, it has been less investigated in the state-of-the-art. Being able to automatically animate the facial expression from a single image would open the door to many new exciting applications in different areas, including the movie industry, photography technologies, fashion, e-commerce business etc. As Generative Adversarial Networks (GANs) have become more successful and more prevalent, a big progress has been made in this task. The most successful architecture is StarGAN, which is able not only to synthesize novel expressions, but also to change other attributes of the face, such as age, hair color or gender. In this project, we will use StarGAN to generate facial expressions corresponding to different emotions. A GAN model in this project is developed to take an image of a person's face and a desired emotion as inputs and as an output one has that person's face with the required emotion applied, while personalized features of the face remain preserved.

II. RELATED WORK

Generating a particular person's face with different facial expressions can be used in a variety of applications, including face recognition [1], [2], face verification [3], [4], emotion prediction, expression database generation, facial expression augmentation and entertainment.

Generative Adversarial Networks (GANs) are a powerful class of generative models based on game theory. A typical GAN optimization scheme consists in simultaneously training a generator network to produce realistic fake samples and a discriminator network trained to distinguish between real and fake data. This idea is embedded by the so-called adversarial loss [5].

DCGANs are generative convolutional networks based off of GANs. Based on work [6], DCGANs seem promising. Using them it is possible to reconstruct the original image with great accuracy. However, it is also shown that the network has not learned to modify the image.

According to the article [5], the GANimation (Anatomically Consistent Facial Animation) proposed additionally controlled generated expressions by Action Units labels, and allowed a continuous expression transformation. The authors introduced

an attention-based generator to promote the robustness of their model for distracting backgrounds and illuminations.

In some related work [7] it was used approach to construct double encoder GAN. Double encoder GAN is used for facial expression synthesis to extract the latent vectors and conditional labels features of the real image.

Recent advances in GANs have shown impressive results for task of facial expression synthesis. The most successful architecture is StarGAN [8], that conditions GANs' generation process with images of a specific domain, namely a set of images of persons sharing the same expression [9]. In our work we have based precisely on this approach.

III. PROPOSED SOLUTION

Our solution is the same as the one proposed in [8]. The StarGAN framework uses a single generator G and a single discriminator D . Generator's role is to produce real images of desired class (in this case desired emotion) conditioned on the input image, while the discriminator is used in the training process of the generator. We denote the output of the generator as $y = G(x, c)$, where y is the generated image, x is the input image and c is the desired class. Out descriminator produces two probability distributions, one over sources and one over classes. We denote those as $D_{\text{src}}(x)$ and $D_{\text{cls}}(x)$.

This is an instance of multi objective optimization because the produced image has to contain the desired emotion, look real as well as similar to the original image. This is reflected in our formulation of the loss function which contains multiple terms, each of which is used to optimize one aspect of the optimization problem. We will consider each term separately and then combine them to formulate the total loss function for the generator as well as for the discriminator.

First term in the loss function is the standard adversarial loss

$$\begin{aligned} \mathcal{L}_{\text{adv}} = & \mathbb{E}_x [\log D_{\text{src}}(x)] + \\ & \mathbb{E}_{x,c} [\log(1 - D_{\text{src}}(G(x, c)))] . \end{aligned} \quad (1)$$

The purpose of this term is to make images generated by the generator look real. The generator tries to minimize this term while the descriminator tries to maximize it. Training the descriminator will in turn make the generator produce work harder to produce images that will be classified as real, therefore the quality of the generated images will increase.

Second term in the loss function is the classification loss. Since we want the generated images to belong to the desired class the classification loss term has to be included in the total

loss function. We can train the discriminator on the real data and thus formulate the loss

$$\mathcal{L}_{\text{cls}}^{\text{r}} = \mathbb{E}_{x,c'} [-\log D_{\text{cls}}(c'|x)]. \quad (2)$$

Since the generator also needs to learn the notion of different classes we can use almost the same loss to train the generator, therefore we formulate the loss

$$\mathcal{L}_{\text{cls}}^{\text{f}} = \mathbb{E}_{x,c} [-\log D_{\text{cls}}(c|G(x,c))]. \quad (3)$$

This expression is almost the same as the one in (2). The only difference is that the real image x from (2) is here replaced by the generated fake image $G(x,c)$. Both the generator and the discriminator want to minimize those losses, although for different reasons. The discriminator wants to get better at classifying, and the generator wants to get better at producing images that will fool the discriminator.

The previous two terms insure that generated images will look real and contain features of the desired class, but there is nothin so far that would stop the generator to learn only one image for each class and completely disregard the x in $G(x,c)$. For this reason we include the third loss term: the reconstruction loss

$$\mathcal{L}_{\text{rec}} = \mathbb{E}_{x,c,c'} [\|x - G(G(x,c), c')\|_1], \quad (4)$$

where c is some class s.t. $c \neq c'$ and c' is the true class of x . This term insures that the generated image will preserve all the features invariant with respect to the class c (in the case of emotions as classes, such features would be e.g. hair color, eye color, nose type...). This is accomplished by applying the generator twice. First time to the image x and with desired class c , and second time to the generated image $G(x,c)$ and with the original image class c' . If the generator changes only features correlated with the class the reconstruction should be easy and the result of cyclical generator application should be almost the same as the original image x , for all the features orthogonal to class will remain unchanged.

The total loss functions for the generator and the discriminator are linear combinations of previously defined terms. The generator should minimize

$$\mathcal{L}_G = \mathcal{L}_{\text{adv}} + \lambda_{\text{cls}} \mathcal{L}_{\text{cls}}^{\text{f}} + \lambda_{\text{rec}} \mathcal{L}_{\text{rec}} \quad (5)$$

and the discriminator should minimize

$$\mathcal{L}_D = -\mathcal{L}_{\text{adv}} + \lambda_{\text{cls}} \mathcal{L}_{\text{cls}}^{\text{r}}, \quad (6)$$

where lambdas are hyperparameters whose variation changes the relative importance of various loss terms, e.g. setting $\lambda_{\text{cls}} = 0$ would probably result in generator producing input images. As in [8] we use $\lambda_{\text{cls}} = 1$ and $\lambda_{\text{rec}} = 10$. Our generator and discriminator architectures are also the same as in [8].

IV. EXPERIMENTAL RESULTS

Dataset: FER2013¹. The data consists of 48×48 pixel grayscale images of faces. The faces have been automatically

cropped so that the face is centered. Images are labeled with 7 different emotions: Angry, Disgust, Fear, Happy, Sad, Surprise and Neutral. We used 4 of them to train our model due to lack of processing power: Angry (3395 images), Happy (7215 images), Sad (4830 images), Surprise (3171 images).

Training was performed using Adam optimizer [10] with $\beta_1 = 0.5$ and $\beta_2 = 0.999$. We train our model with learning rate 0.0001 for 400 epochs. Batch size is set to 8. Training execution time was roughly 1 day on NVIDIA Tesla T4 GPU.

Gradual improvement of desired facial expressions can be seen throughout epochs. After 100 epochs, there are only slight changes which are prevalent in the mouth area. Visible lines, which remind of moustache, appear around the mouth since it is where facial expressions most differ. By the 400th epoch, modifications are more obvious and the desired emotions more easily recognizable. Changes involve several facial areas, including eyebrows, eyes, etc.

V. CONCLUSION

Conclusion...

REFERENCES

- [1] Cao, Q., Shen, L., Xie, W., Parkhi, O.M., & Zisserman, A. (2018). Vggface2: A dataset for recognising faces across pose and age. In 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018) (pp. 67–74). IEEE.
- [2] Parkhi et al. 2015 - Parkhi, O. M., Vedaldi, A., & Zisserman, A. (2015). Deep face recognition. In BMVC (Vol. 1, p. 6).
- [3] Sun et al. 2014; - Sun, Y., Chen, Y., Wang, X., & Tang, X. (2014). Deep learning face representation by joint identification–verification. In Advances in neural information processing systems (pp. 1988–1996).
- [4] Taigman et al. 2014 - Taigman, Y., Yang, M., Ranzato, M., & Wolf, L. (2014). Deepface: Closing the gap to human-level performance in face verification. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1701–1708).
- [5] Pumarola, A., Agudo, A., Martinez, A.M. et al. GANimation: One-Shot Anatomically Consistent Facial Animation. Int J Comput Vis 128, 698–713 (2020).
- [6] Sut, J. Generating Facial Expressions
- [7] Chen, Mingyi & Li, Changchun & Li, Ke & Zhang, Han & He, Xuanji. (2018). Double Encoder Conditional GAN for Facial Expression Synthesis. 9286-9291. 10.23919/ChiCC.2018.8483579.
- [8] Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. CVPR (2018)
- [9] Pumarola, A., Agudo, A., Martinez, A.A., Sanfelix, A., Moreno-Noguer, F.: Ganimation: Anatomically-aware facial animation from a single image. In: ECCV (2018)
- [10] Kingma, Diederik & Ba, Jimmy. (2014). Adam: A Method for Stochastic Optimization. International Conference on Learning Representations.

¹<https://www.kaggle.com/msambare/fer2013?select=train>