

Image Inpainting with Generative Models

Daniyar Nariman, Sultanov Azat, Nazyrov Amir

Innopolis University

Course: Computer Vision 2019

Abstract—Recent advances in Deep Learning have shown an exciting promise in filling the gaps in the natural images. In computer vision this kind of tasks called image inpainting, which aims to synthesize alternative contents in missing regions based on the context. Many approaches have been proposed on this kind of task, that confirms about its relevance in the current days, but it still remains a challenging problem. In this project we will briefly explain the current approaches of Image Inpainting task and show the working prototype of one of the methods.

I. INTRODUCTION

Image inpainting is a task of synthesizing visually realistic and semantically correct contents in missing regions.[1] There are mainly two broad approaches: classical patch-based ones and convolutional networks based deep generative models. Patch based approach can synthesize plausible stationary textures but typically does not work well in cases where semantics are needed. The CNN based approach can utilize semantics learned from large scale datasets. In this paper we will focus on the second approach as it seems more promising. Also we describe new EdgeConnect approach which used adversarial edge learning [4]

II. APPROACHES

In this section we will describe some of the approaches, which was developed for the past few years, but before going deeply on the methods of Image Inpainting, we will shortly explain abstract workflow.

- Feed the generative network with an input image with holes than need to be filled. These patches can be considered a hyperparameter required by the network since the generator has no way of discerning what actually needs to be filled in. To make network understand what part need to be filled, we need a separate

layer mask that contains pixel information for the missing data.

- The generator will produce the entirely synthetic image generated from scratch.
- The layer mask allows us to discard those portions that was already presented in the incomplete image, since we did not interested in those parts of the image.
- The new generated image is then combined with the incomplete image that we sent to the network at the beginning, so that we can get our new image that generated pixels at the place where the hole was.
- This generated image is then sent to discriminator network. This network tries to ensure that the image does not look obviously fake.

A. Context Encoder [2] [3]

It is a convolutional neural network trained to generate the contents of an arbitrary image region conditioned on its surroundings.

1) *Generator*: The overall architecture is a simple encoder-decoder pipeline. The encoder was to be derived from the *AlexNet* architecture and composed of five convolutional layers. The output of the encoder is connected to the input of the decoder by a channel-wise fully-connected layer ($6*6*256=9126$ activations layers), in order to propagate the information within each feature map. The decoder follows with five up-convolutional layers that augment the dimensions of the latent representation in-between the encoder and the decoder.

2) *Discriminator*: Five layer convolutional network that takes as input the context either completed with the output of the encoder-decoder pipeline (which is the generated missing area) or the ground truth missing area.

3) *Loss function*: The overall joint loss function is shown below.

$$\mathcal{L} = \lambda_{rec}\mathcal{L}_{rec} + \lambda_{adv}\mathcal{L}_{adv}$$

where the first part is *reconstruction loss* which is equal to

$$\mathcal{L}_{rec}(x) = ||\hat{M} * (x - F((1 - \hat{M}) * x))||_2^2$$

where x designate the input image and \hat{M} the binary mask such that $\hat{M} * x$ (Hadamard product) selects the missing region. The reconstruction loss is the L2 distance between the content generated by the network and the original region of the image that has been masked to the network, and the second part is represented as,

$$\mathcal{L}_{adv}(F, X) = E_{X \in \mathcal{X}}[\log D(X)] \\ + \log(1 - D(F((1 - \hat{M}) * x)))$$

F here is the function of our proposed context encoder and λ is hyperparameter.

The main advantage of this approach is that, because of encoder-decoder based architecture, it learns a representation that captures not just appearance but also the semantics of visual structures, that gives the model more power in the case of image inpainting task.

B. EdgeConnect [4]

Idea of this approach is simple: model tries to get edges of the mask image, restore edges and then generate images by this edge map. EdgeConnect consist of two stages each of which use generator/discriminator: edge generator and image completion network. The edge generator hallucinates edges of the missing region of the image, and the image completion network fills in the missing regions using hallucinated edges as a priori.

1) *Edge generator*: Edge generator takes grayscale image as input, its edge map, and image mask as a precondition and by this it can predict the edge map for the masked region: $C_{pred} = G_1(\tilde{I}_{gray}, \tilde{G}_{gt}, M)$. They use C_{gt} and C_{pred} as inputs of the discriminator that predicts whether or not an edge map is real.

Network was trained using adversarial loss and feature-matching loss:

$$\min_G \max_D \mathcal{L}_G = \min_G (\lambda_{adv,1} \max_D (\mathcal{L}_{adv}, 1) + \lambda_{FM} \mathcal{L}_{FM})$$

where:

$$\mathcal{L}_{adv,1} = E_{(C_{gt}, I_{gray})} [\log D(C_{gt}, I_{gray})] \\ + E_{I_{gray}} \log [1 - D(C_{pred}, I_{gray})]$$

$$\mathcal{L}_{FM} = E \left[\sum_{i=1}^L \frac{1}{N_i} \left\| D_1^{(i)}(C_{gt}) - D_1^{(i)}(C_{pred}) \right\|_1 \right]$$

2) *Image Completion network*: Input to the image completion network is incomplete color image \tilde{I}_{gt} and composite edge map C_{comp} . Edge map then is combined by the background edges of the initial image and generated edges from the previous step. The generator returns the color image with missing region filled in:

$$I_{pred} = G(\tilde{I}_{gt}, C_{comp})$$

Discriminator predicts the realistic of the output image.

Loss function consist of several loss functions: l_1 loss, adversarial loss which is similar adversarial loss from previous step, perceptual loss and style:

$$\mathcal{L}_{adv,2} = E_{(I_{gt}, C_{comp})} [\log D_2(I_{gt}, C_{comp})] \\ + E_{C_{comp}} \log [1 - D_2(I_{pred}, C_{comp})]$$

$$\mathcal{L}_{perc} = E \left[\sum_i \frac{1}{N_i} \left\| \phi_i(I_{gt}) - \phi_i(I_{pred}) \right\|_1 \right]$$

$$\mathcal{L}_{style} = E_j \left[\left\| G_j^\phi(\tilde{I}_{pred}) - G_j^\phi(\tilde{I}_{gt}) \right\|_1 \right]$$

So overall loss function is:

$$\mathcal{L}_{G_2} = \lambda_{\ell_1} \mathcal{L}_{\ell_1} + \lambda_{adv,2} \mathcal{L}_{adv,2} + \lambda_p \mathcal{L}_{perc} + \lambda_s \mathcal{L}_{style}$$

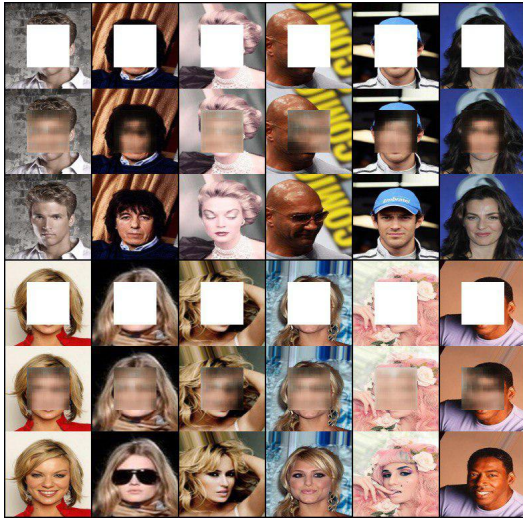


Fig. 1. Each image sample has 3 states: cropped, generated from our Context encoder, original

3) Generator and Discriminator architecture:

The structure of the generators and discriminators for both GANs the same. The Generator use in-network downsampling to reduce the spatial extent of feature maps followed by in-network upsampling to produce the final output image. While for discriminators, they use a 70×70 Patch-GAN architecture, which determines whether or not overlapping image patches of size 70×70 are real.

The main advantage of this approach is that GAN restore edges. Since the mask is not on the whole image, but only on a part, some edges have a continuation on the other side, so edge generator learn to continue edges to the other side. But also it has disadvantage: in this type of architecture we have 2 different GANs, so time on the learning increase.

C. Another Approach

another approach to be added

III. EXPERIMENTS AND EVALUATION

A. Implementation

As the part of the research, before training on our dataset, we conducted an experiment on Context encoder GAN, by training it on 240 images. The results are represented in Fig. 1.

The implementation of Context encoder GAN you can find in the following link:

<https://github.com/nariman9119/image-inpainting-research>

The next checkpoint for our team, is to train this model on generating the mustaches on the cropped images.

B. Dataset

The training that we planning to use is a Large-scale CelebFaces Attributes (CelebA) Dataset with the attribute "Mustache". The number of images in this dataset 8417.

IV. CONCLUSION

TODO

REFERENCES

- [1] Tianyuan Zhang, 2018. Survey on Image Inpainting.
- [2] Deepak Pathak, Philipp Krahenb, Jeff Donahue, Trevor Darrell, Alexei A. Efros, 2016. Context Encoders: Feature Learning by Inpainting
- [3] Quentin LEROY and Bastien PONCHON 2018. Understanding Context encoders: Feature Learning by Inpainting.
- [4] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z. Qureshi, Mehran Ebrahimi, 2019. EdgeConnect: Generative Image Inpainting with Adversarial Edge Learning