

CSE-578 Computer Vision

Project Report Context Encoders: Feature Learning by Inpainting

March 8 2020

1 Abstract

In this paper, the authors propose a novel unsupervised feature learning method using inpainting. Inpainting is the task of filling a potentially large region of an image (one fourth for example). The method is Context Encoder, a convolutional neural network trained for inpainting. Trained to create plausible content for the missing part of an image based on the surroundings, the network learns a representation of the input that captures the appearances and semantics of visual structures. The features learned by the network are shown to be effective for other tasks than inpainting that usually requires supervised learning for example segmentation, detection and classification.

1.1 Introduction

Convolutional Neural Networks (CNN) were originally found to be effective for the classification task under a supervised setting. A tremendous amount of later works further proved that they also excel for other computer vision tasks. In the present work the authors explore CNN trained in an unsupervised manner. The network is asked to completely restore a large part of the input image that is missing (for example a $64 * 64$ square in the center of an $128 * 128$ image). This enforces the network to holistically understand the semantics of the scene, learns high-level features over the whole spatial extent of the image in order to provide content for the missing part that semantically complies with the context of the scene. This is achieved thanks to an Auto-Encoder-like architecture. The network is composed of an **Encoder** that transforms the input image into a fixed small number of features followed by a **Decoder** that creates out of these features the content of the missing part. The parameters are learned with a reconstruction loss as well as an adversarial loss.

1.2 Context Encoder-Decoder: Network Architecture

The original architecture of Context encoders proposed by the authors was to use an autoencoder-like structure to code the context of an input image as a $6 * 6 * 256$ feature map (9126 hidden features). 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 (at the interface of the encoder and the decoder).

1.2.1 Loss Function

The choice of the loss function is strongly related with the task at hand namely inpainting or put it differently filling a hole missing in the input image. The paper initially uses a regular $L2$ norm based loss function, but the resulting patches were blurry and were not sharp. To state in detail, we have x as our input image and \hat{M} the binary mask, such that $\hat{M} * x$ (Hadamard product) selects the missing region. The formula is as follows:

$$L_{recx} = || \hat{M} * (x - F(1 - \hat{M}) * x) ||^2$$

Note that the mask \hat{M} can be of any size and shape, covering up to $1/4$ of the image. This pixel-wise distance does not encourage the network to produce details with high precision but rather create content that overall captures the structure of the original region. To overcome the blurry result favored by the reconstruction loss the authors propose to incorporate an adversarial loss to the total loss.

1.3 GAN Framework : Adversarial Loss

The adversarial loss is based on Generative Adversarial Networks. GANs were introduced by [6] as deep generative models able to learn to represent an estimate, either explicitly or implicitly, of some distribution $p_{data}(X)$ on some space X , given a training set composed of samples x drawn from this distribution. As explained in [5], this is done through a two player game between two parameterized and differentiable (with regards to their inputs and their parameters) functions: a generator G and a discriminator D .

- The **Generator $G : Z \rightarrow X$** : takes as input some z drawn from a prior distribution $p(z)$ on a space Z and tries to produce an output $G(z) \in X$ drawn from a probability distribution p_{model} such that that p_{model} is close from p_{data} , that is to say such that G_z seems to have been drawn from the distribution $p_{data}(X)$.
- The **Discriminator $D : X \rightarrow [0, 1]$** : on the other hand takes as input an element $x \in X$ and outputs the probability $D(x)$ that x comes from the training set, that is to say, has been drawn from the true distribution $p_{data}(X)$, rather than having been generated by G , and therefore drawn from the distribution $p_{model}(X)$. This translates by the maximization of the loss function L_{adv} given by the following log likelihood (or the minimization of the corresponding negative log likelihood):

$$\begin{aligned} \mathcal{L}_{adv}(G, D, Z, X) = & E_{x \sim p_{data}(X)} (\log D(x)) \\ & + E_{z \sim p_Z} (\log(1 - D(G(z)))) \end{aligned}$$

The generator tries to fool the discriminator, and therefore tries to minimize this very value $L_{adv}(G, D, Z, X)$. The solution of this mini- max game (or zero-sum game, as the generator and the discriminator try to minimize two opposite cost functions) is given by $\min_G \max_D L_{adv}(G, D, Z, X)$.

1.4 Adversarial Loss

In the case of our inpainting task, the target domain X is the space of complete(d) images, with p_{data} represented by the set of training images. The source domain Z is the space of context images (image with a missing masked region). The encoder-decoder pipeline F aims at generating a complete image out of an incomplete one and therefore can be seen as a generator. The authors hence introduce a corresponding discriminator network D that aims at predicting whether the input fed is a real sample or is a sample coming from the distribution of the content generated by the encoder-decoder pipeline. The discriminator is a 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. The adversarial loss is then :

$$\begin{aligned} \mathcal{L}_{adv}(F, D) = & \mathbb{E}_{X \in \mathcal{X}} [\log D(x)] \\ & + \log(1 - D(F((1 - \hat{M}) * x))) \end{aligned}$$

1.4.1 Joint Loss

The joint loss that the encoder-decoder generator tries to minimize is a weighted combination of the reconstruction loss and the adversarial loss .

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

Here, the λ_1 value is set to 0.1 and λ_2 value is set to 0.9. Note that unlike the base GAN framework, the loss that the generator of the context encoder tries to minimize is not strictly equal to the opposite.

2 Model Diagram : Fixed Size Patch

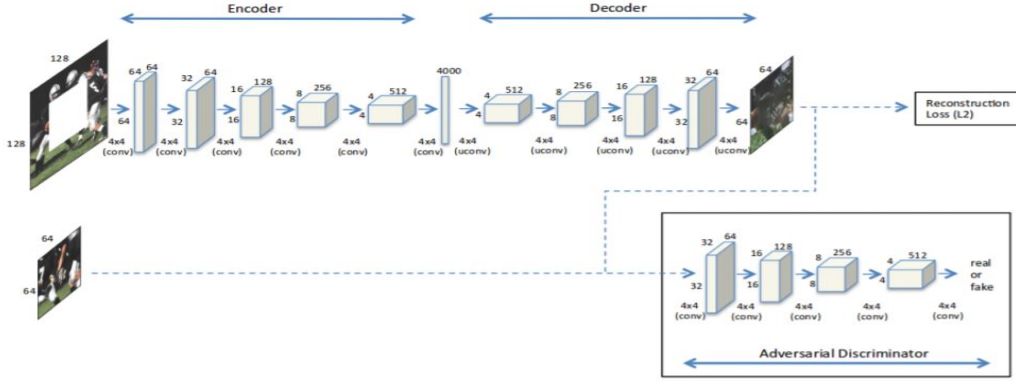


Figure 1. Context Encoder trained with joint reconstruction and adversarial loss for **semantic inpainting**. This illustration is shown for center region dropout. Similar architecture holds for arbitrary region dropout as well

Figure 1: Model Architecture

3 Our Code Results

For producing the following results , we have taken up CIFAR-10 dataset which contains 10 classes ,with 6000 images in each category per class ,we have specifically chosen CAT DOG class data.The entire training set consists of 50000 images and test set consists 10000 images. The network has undergone 30000 epochs while training.

3.1 Encoder-Decoder Model

```

1 def context_encoder(input_image):
2     # Encoder
3     network = conv_2d(input_image, 32, 4, strides=2, activation='relu')
4     network = conv_2d(network, 16, 4, strides=2, activation='relu')
5     network = conv_2d(network, 8, 4, strides=2, activation='relu')
6     # network = conv_2d(network, 512, 4, strides=2, activation='relu')
7
8     network = max_pool_2d(network, 6, strides=2)
9
10    # Channel wise connection : RGB
11    network0 = fully_connected(network[:, :, :, 0], 16) #Channel ----> R
12    network1 = fully_connected(network[:, :, :, 1], 16) #Channel ----> G
13    network2 = fully_connected(network[:, :, :, 2], 16) #Channel ----> B
14
15
16    #Decoder
17    network = tf.reshape(tf.concat([network0, network1, network2], axis=1), (-1, 4, 4, 3))
18
19    network = conv_2d_transpose(network, 12, 4, [4,4], strides=1, activation='relu')
20    network = conv_2d_transpose(network, 12, 4, [4,4], strides=1, activation='relu')
21    network = conv_2d_transpose(network, 6, 4, [8,8], strides=2, activation='relu')
22    network = conv_2d_transpose(network, 3, 4, [16,16], strides=2, activation='relu')
23    # network = conv_2d_transpose(network, 3, 4, [64,64], strides=2, activation='relu')
24    return network
25
26 def discriminator(input_image):
27     network = conv_2d(input_image, 32, 4, activation='relu')
28     network = conv_2d(network, 64, 4, strides=2, activation='relu')
29     network = conv_2d(network, 128, 4, strides=2, activation='relu')
30     # network = conv_2d(network, 512, 4, strides=2, activation='relu')
31     network = fully_connected(network, 1)
32     return network

```

3.2 Qualitative Results

We present the following images at different instances of training epochs, arranged in the order of decreasing training epoch number .

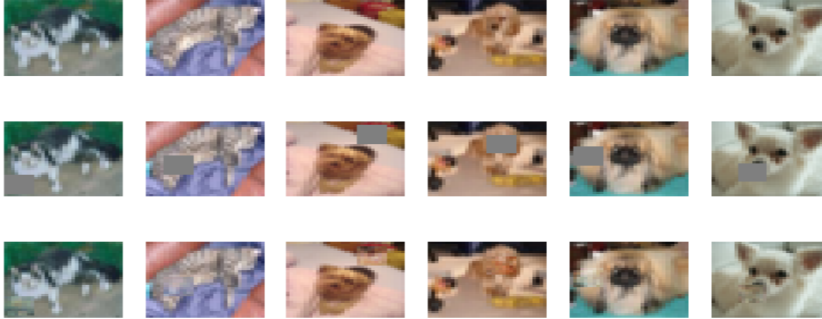


Figure 2: Iteration : 29900

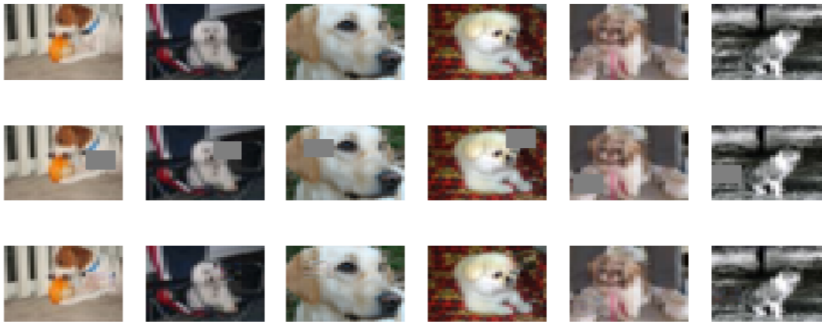


Figure 3: Iteration : 29300

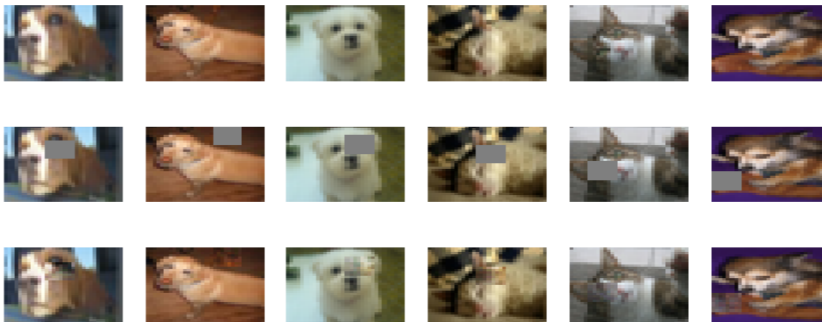


Figure 4: Iteration :20900

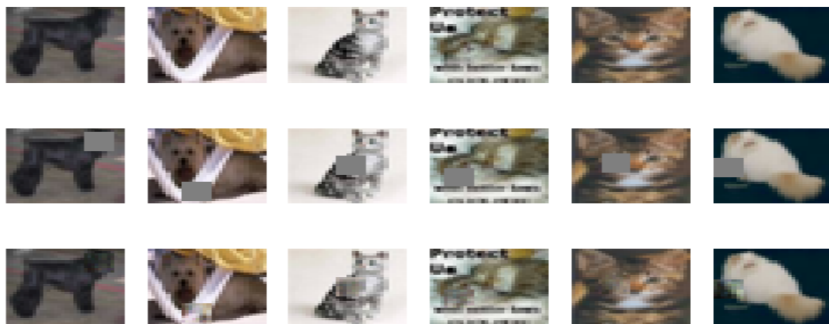


Figure 5: Iteration :15500

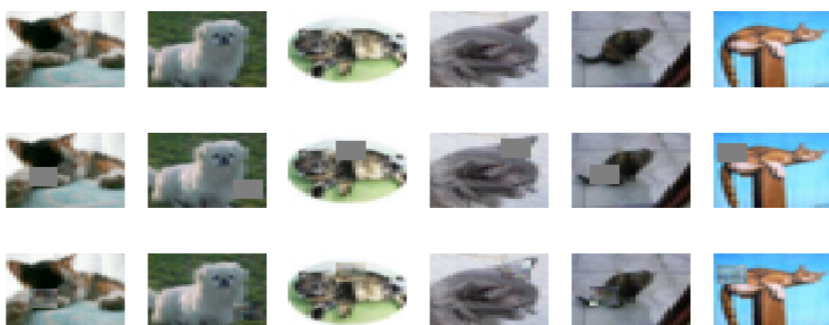


Figure 6: Iteration :10100

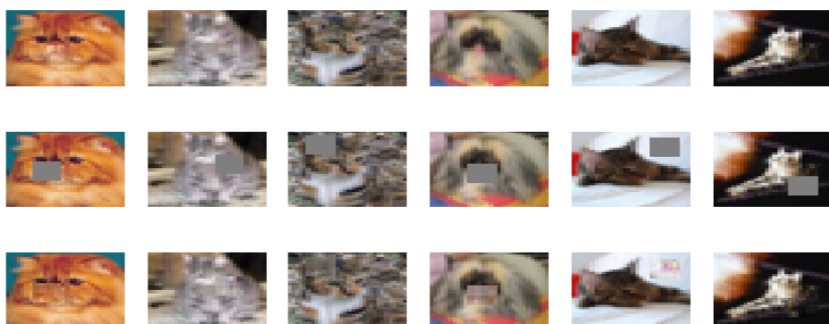


Figure 7: Iteration :5300

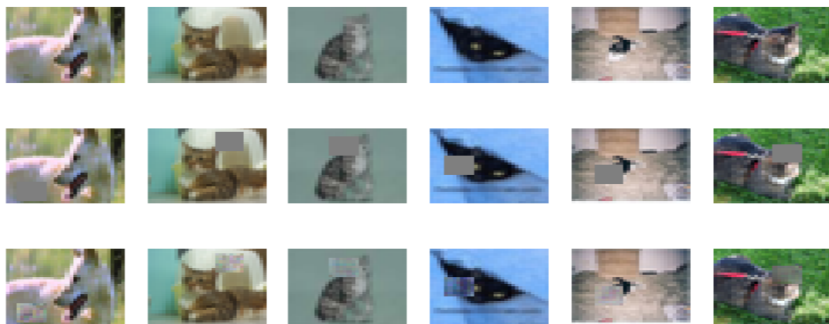


Figure 8: Iteration :1700

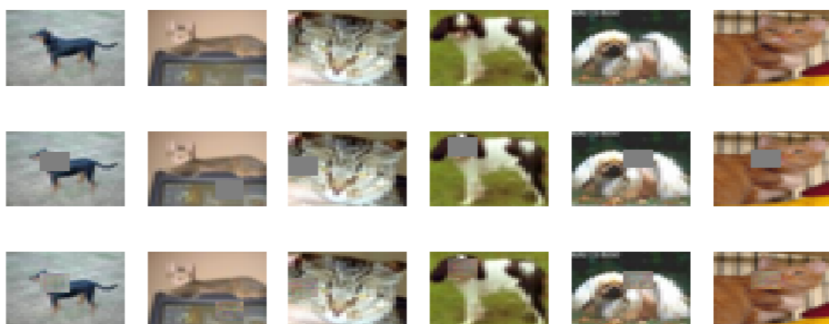


Figure 9: Iteration :800

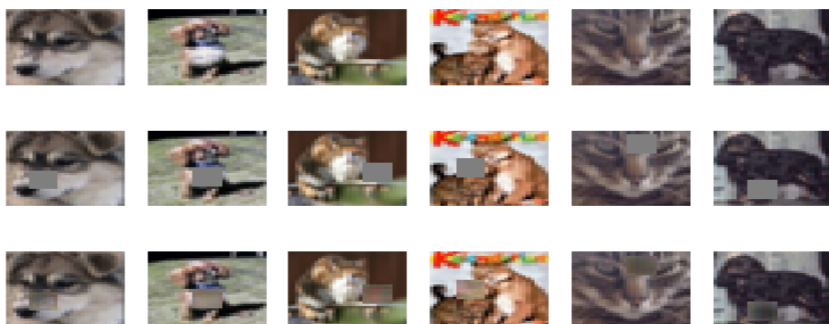


Figure 10: Iteration :200

Note: Link to all epoch training images can be obtained here :

https://drive.google.com/drive/folders/13_RU8N9UIMoVW0tpwsTXkooqdbsqy_so?usp=sharing

Link to the Colab Notebook :

<https://colab.research.google.com/drive/1HGLopRF0NORxHGfHBy7MW2K7hQIKHwa3#scrollTo=1RNETl0mM57J>

4 Remaining Work

We can see above that we have scaled down the entire network to quarter size for faster training and only 2 classes in CIFAR 10 was chosen . Our goal is to use diverse dataset with $128 * 128$ sized images. Also, our next part of the project will include removing random and arbitrary shape objects.

5 References

- <https://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>
- <https://people.eecs.berkeley.edu/~pathak/papers/cvpr16.pdf>
- <http://ufldl.stanford.edu/tutorial/unsupervised/Autoencoders/>

6 Team

- S.Dhawal - 2019201089
- Sravya.S - 2019702008
- Niharika.V - 2018122008