

National Chiao Tung University

Spring 2019

Deep Learning

Instructor: Jen-Tsung Chien

Deep Learning Final Project Report

Neural Style Transfer to Make Photo Younger

Group 13 members

Alfons Hwu (Group leader)

0416324

Pin-Jen Hunag

0416033

Chia-Yu Sun

0416045

Min-Xue Yang

0416022

Dept. of Computer Science

Composed with L^AT_EXon Overleaf

Deep Learning Final Project Report

June 20, 2019

1 Introduction

Here we introduce an artificial system based on a Convolutional Neural Network that creates artistic images of high perceptual quality. The system uses neural representations to separate and recombine content and style of arbitrary images, providing a neural algorithm for the creation of artistic images. In this project, we hope to "translate" the image from "young face" to "old face" and vice versa by using the "Neuron Style Transfer" mentioned in <https://arxiv.org/abs/1508.06576> with vgg network and some self-designed model.

2 Model Architecture

The architecture is simple, we have three images, S represents the style image, C represents the content image and G represent the synthesized image to be generated.

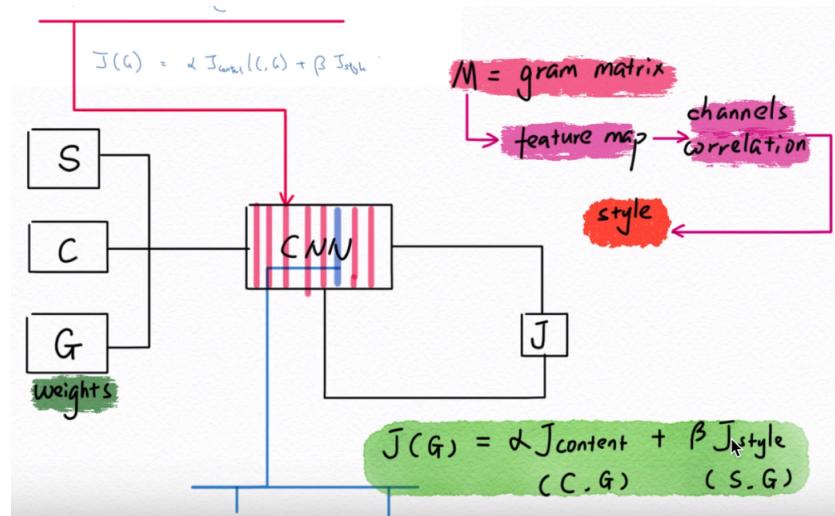


Figure 1: Model architecture

The CNN model used is the pretrained VGG-19 network, which is a convolutional neural network

that is trained on more than a million images from the ImageNet database. The network is 19 layers deep and trained on millions of images. Because of which it is able to detect high-level features in an image.

During the convolution procedure, features in images can be filtered out and they are called "feature maps". We may calculate the content loss and style loss respectively with p being the content image, a being the art style image, x being the generated, synthesized image.

Hence, overall loss to be

$$L_{total}(p, a, x) = \alpha L_{content}(p, x) + \beta L_{style}(a, x)$$

3 Understand the math behind the loss

Reference and revised from <https://arxiv.org/abs/1508.06576> page 10 to 12

3.1 The loss of content

A layer with N_l distinct filters has N_l feature maps each of size M_l , where M_l is the height times the width of the feature map (magnitude).

So the responses in a layer l can be stored in a matrix $F_l \in R^{N_l \times M_l}$ where F_{ij}^l is the activation of the i^{th} filter at position j in layer l and the same is true for P_{ij}^l

Let p and x be the original image and the image that is generated. P^l, F^l being their respective feature representation in layer l after extracted by convolutional filter.

We may thus use the mean square error to compare the original content image and generated image pixel-wisely.

$$L_{content}(p, x, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2$$

3.2 The loss of style

On top of the CNN responses in each layer of the network we built a style representation that computes the correlations between the different filter responses, where the expectation is taken over the spatial extend of the input image. These feature correlations are given by the Gram matrix $G_l \in R^{N_l \times N_l}$, where G_{ij}^l is the inner product between the vectorised feature map i and j in layer l :

$$G_{ij}^l = \sum_k (F_{ik}^l - F_{jk}^l)$$

Finally we use gradient descent to minimize the difference between original style image a and generated image x and A^l and G^l being their respective style representations in layer l

Total loss of style will be

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

Why does it work?, especially in the style loss.

First, try to optimize the content difference is quite intuitive.

Yet concerning the style loss,in order to capture the style of an image we would calculate how “correlated” these filters are to each other meaning how similar are these feature maps.

Here the **Gram matrix** helps solve this problem, since it is defined as a set of vectors in an inner product space, which is the **Hermitian matrix** of inner products, whose entries are given by $G_{ij} = \langle V_i, V_j \rangle$. If the dot-product across the activation of two filters is large then two channels are said to be correlated and if it is small then the images are un-correlated.

如何用简单易懂的例子解释格拉姆矩阵/Gram matrix?

在看CNN style transfer遇到的问题，数学不好只好上来问问大家啦

关注问题 写回答 邀请回答 添加评论 分享 ...

10 个回答 默认排序

 **90后后生**
没有一天不懵B的研究僧

75 人赞同了该回答

Gram Matrix实际上可看做是feature之间的偏心协方差矩阵（即没有减去均值的协方差矩阵），在feature map中，每一个数字都来自于一个特定滤波器在特定位置的卷积，因此每个数字就代表一个特征的强度，而Gram计算的实际上是两两特征之间的相关性，哪两个特征是同时出现的，哪两个是此消彼长的等等，同时，Gram的对角线元素，还体现了每个特征在图像中出现的量，因此，Gram有助于把握整个图像的大体风格。有了表示风格的Gram Matrix，要度量两个图像风格的差异，只需比较他们Gram Matrix的差异即可。

Figure 2: Zhihu reply for understanding the gram matrix

Based on the reply from <https://www.zhihu.com/question/49805962>

$G_{ij}^l = \sum_k (F_{ik}^l - F_{jk}^l)$ will go through each combination of inner-producted feature maps, either in the original image or generated image, consequently all the vital features of images are represented. Cost function between style and generated image is the **square of difference between the Gram Matrix of the style Image with the Gram Matrix of generated Image**. $E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$ Finally we will derive the total loss function $L_{total}(p, a, x) = \alpha L_{content}(p, x) + \beta L_{style}(a, x)$ By adjusting the weight of α and β , we may either emphasize the style or content for the generated image.

4 Dataset Description

We have collected some of the young face as the style image and old face as content image hoping to rejuvenate the old face.



Figure 3: Dataset images

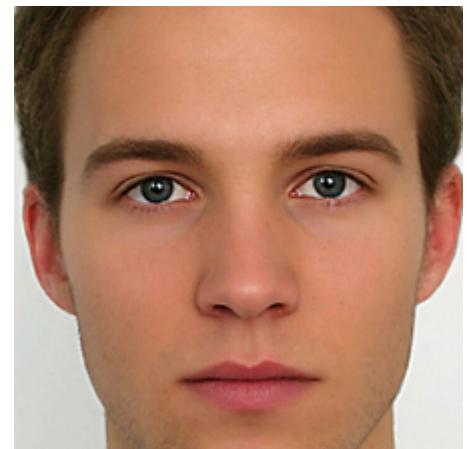
5 Dataset Preprocessing

6 Result Evaluation

Following are some schemes that we evaluate as overall optimal in our project for report and demo.



(a) Old face image(content)



(b) Young face image(style)

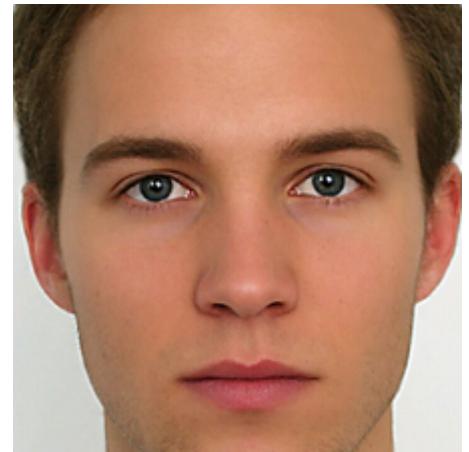


(c) Neuron style transferred image

Figure 4: Style transfer (make photo younger) scheme 1



(a) Old face image(content)



(b) Young face image(style)



(c) Neuron style transferred image

Figure 5: Style transfer (make photo younger) scheme 2

7 Group Work Distribution

0416324: NN architecture, testing, shell scripting to automatically organize and train data / IO, writing this LaTeX report.

0416033: Organize data, checking output result and feedback to 0416324 and 0416045 to revise.

0416045: NN architecture, testing and ppt.

0416022: Make demo video, ppt

8 Reference Data

Neuron Style Algorithm: <https://arxiv.org/abs/1508.06576>

Neuron Style Superposition(famous paper by Li Fei-Fei from Stanford University) <https://arxiv.org/abs/1508.06576>

Zhihu reply <https://www.zhihu.com/question/49805962>

Medium article 1 <https://hackernoon.com/how-do-neural-style-transfers-work-7bedaee0559a>

Medium article 2 <https://towardsdatascience.com/style-transfer-styling-images-with-convolutional-neural-networks-10f3a2a2a2d>

Bilibili video tutorial <https://www.bilibili.com/video/av22581633/?p=41>

Quora reply <https://www.quora.com/What-is-Gram-matrix-What-is-its-significance-in-linear-algebra>