

Extracted the style from the image with Deep Neural Network

吴秉哲

2015 年 12 月 9 日

摘要

在视觉艺术的领域，尤其是在绘画方面，画家已经能通过自己的经验和
技术进行创作出一副独特的复杂的艺术作品，这些作品的内容和风格的联
系十分复杂。在计算机视觉领域，由于卷积神经网络的兴起，一些传统的任
务比如物体定位和物体识别已经取得了很大的进步。尤其是物体识别已
经有算法在大规模图片识别 (imageNet 图片数据集) 的任务上达到了类人
的表现 [4]，其他方面，由于 GPU 通用计算 [1] 的发展，人们已经有能力进
行深度人工神经网络的训练 (2014 年 ImageNet 挑战赛的冠军使用的网络
模型已经到了 19 层)。在这次的期末项目中，我会使用深层的卷积神经网
络来提取一副图像的物体和另一副绘画的艺术风格，并将它们结合在一起，
形成一副全新的图片。项目中使用的算法也提供一个途径帮助计算机理解
人类是如何完成复杂的艺术创作的。

1 Introduction

在这次的项目中，我们主要使用一种特殊的人工神经网络，卷积神经网络。卷积神经网络可以理解为一个特征提取的网络，区别与传统的计算机视觉的方法，卷积神经网络会自动完成对图片特征的提取，在训练一个多层的卷积神经网络的时候，随着层数的增加，图片中物体的轮廓会越来越清晰，但是一些图片的细节会损失。根据卷积神经网络这个特点，我们在提取一张图片的具体内容（就是图片里面的一些物体，比如建筑，动物等）的时候，我们只需要根据每一个卷积层得到的特征图（即上一层图片进过这层卷积得到的图片），对特征图进行重建，我们就可以提取图片中内容的信息。在这次的期末项目中，我们主要利用网络中深层所得到的特征图进行重建得到图片的内容的信息。
To obtain a representation of the style of an input image ,we use a feature space originally designed to capture texture information.This feature space we used is built on top of the filter responses in each layer of the network .It consists of the correlations between the different filter responses over the spatial

extent of the feature maps.(See Implementation for details).By including the feature correlations of multiple layers,we obtain a multi-scale representation of the input image ,which captures its style information but not the details of content and its arrangement.

Again,we can visualise the information captured by these style feature spaces built on different layers of the network by constructing an image that matches the style representation of a given input image .Indeed reconstruction from the style features from the style features produce texturised versions of the input image that capture its general appearance in terms of color and localised structures.Moreover,the size and complexity of local image structures from the input image increases along the hierarchy,a result that can be explained by the increasing receptive field sizes and feature complexity.We refer to this multi-scale representation as style representation.

In this project ,we can manipulate both representation independently to produce new meaningful image.To demonstrate this work,we generate image mix the content and style representation from two different source images.In particular.

The images are synthesised by finding an image that simultaneously matches the content representation of the photograph and the style representation of the respective piece of art.While the global arrangement of the original photograph is preserved,the colors and local structures that compose the global scenery are provided by the artwork.Effectively,this renders the photograph in the style of the artwork,such that the appearance of the synthesised image resembles the work of art,even though it shows the same content as photograph.

Of course,image content and style cannot be completely disentangled.When synthesising an image that combine the content of one image with the style of another image,there usually does not exist an image that perfectly matches both constraints at the same time.However,the loss function we minimise during image synthesis contains two terms for content and style respectively,that are well separated.We can therefore smoothly regulate the emphasis on either reconstructing the content or the style.

In my final project ,we set up an artificial system that achieves a separation of image content from style,thus allowing to recast the content of one image in the style of any other image. We demonstrate this by creating new,artistic images that combine the style of several well-known paintings with the content of an arbitrarily chosen photograph that I take.In particular,we derive the neural representations for the content and style of an image from the feature responses of high-performing Deep Neural Networks trained on object recognition.

As outlined above,

The rest of the report is organized as follows :Section 2 provides a background for CNN and VGG model,and also described the datasets for training model we used.Section 3 presents our described the method we used .Section 4 describes implementation details.Section5 shows our experiment result.

2 BACKGROUND

2.1 DATASETS

ImageNet is a dataset of over 15 million labeled high-resolution images belong to roughly 22,000 categories.The images were collected from the web and labeled by human labelers using Amazon's Mechanical Turk crowd-sourcing tool.Starting in 2010 ,as part of the Pascal Visual Object Challenge,an annual competition called the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) has been held.ILSVRC uses a subset of ImageNet with roughly 1000 images in each of 1000 categories.In all,there are roughly 1.2 million training images,50,000 validation images,and 150,000 testing images.

ImageNet consists of variable-resolution images,while our system requires a constant input dimensionality .In our training process,we down-sampled the images to a fixed resolution of 256×256 .Given a image , we first rescaled the image such that the shorter side was of length 256,and then cropped out the central 256×256 patch from the resulting image.We did not pre-process the images in any other way,except for subtracting the mean activity over the training set from each pixel. Then we use these picture for training a VGG19(see it in next subsection) model by Caffe[3](A deeplearning framework).

2.2 CNN Basics

Convolutional neural network (CNN) is first inspired by research in neuroscience.After over twenty years of convolution,CNN has been gaining more and more distinction in research fields,such as computer vision,pattern recognition ,NLP.As a classical supervised learning algorithm,CNN employs a feedforward process for recognition and backward path for training.

A typical CNN is composed of two components : a feature extractor and a classifier . The feature extractor is used to filter input images into feature maps that represent various features of the image.These features may include corners,lines,circular arch,etc.,which are relatively invariant to position shifting or distortions.The output of the feature extractor is a low-dimensional vector

containing these features. This vector is then fed into the classifier, which is usually based on traditional artificial neural networks. The purpose of this classifier is to decide the likelihood of categories that the input image might belong to.

A typical CNN is composed of multiple computation layers. For example, the feature extractor may consist of several convolutional layers and optional sub-sampling layers. Figure 1 illustrates the computation of a convolutional layer. The convolutional layer receives N feature maps as input. Each input feature map is convolved by a shifting window of size S , which is normally smaller than K . A total of M output feature maps will form the set of input feature maps for the next convolutional layer.

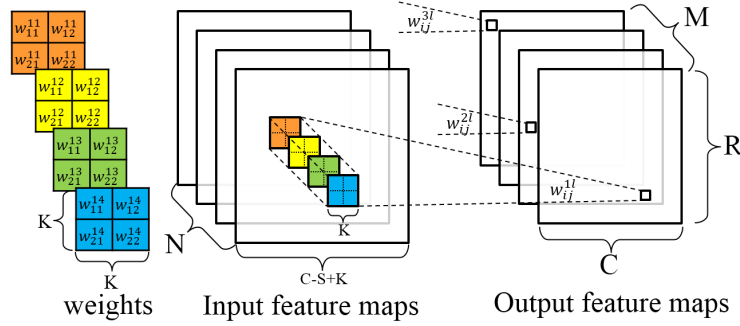


图 1: Graph of a convolutional layer

2.3 VGG Model

2.3.1 ARCHITECTURE

VGG model is a convolutional neural network that rivals human performance on a common visual object recognition benchmark task and was introduced and extensively described in [5].

In this project, we used VGG [5] models for training. During training, the input to this ConvNets is a fixed-size 224×224 RGB image. The only pre-processing we do is subtracting the mean RGB value, computed on the training set, from each pixel. The image is passed through a stack of convolution layers, where we use filters with a very small receptive field: 3×3 (which is the smallest size to capture the notion of left/right, up/down, center). In one of the configurations we also utilise 1×1 convolution filters, which can be seen as a linear transformation of the input channels (followed by non-linearity). The convolutional stride is fixed to 1 pixel; the spatial padding of conv. layer input is such that the spatial resolution is preserved after convolution, i.e. the padding is

1 pixel for 3×3 conv.layers. Spatial pooling is carried out by five max-pooling layers, which follow some of the conv.layers (not all the conv.layers are followed by max-pooling). Max-pooling is performed over a 2×2 pixel window, with stride 2.

A stack of convolutional layers (which has a different depth in different architectures) is followed by three Fully-Connected (FC) layers. However in our project, the full-connected layers were removed.

All hidden layers are equipped with the rectification (ReLU) non-linearity. We note that none of our networks (except for one) contain Local Response Normalisation (LRN) normalisation.

The Convnet we used in the project are shown in Figure 2

3 Methods

The results presented in our project were generated on the basis of the VGG-Model (see details in section 2). We used the feature space provided by the 16 convolutional and 5 pooling layers of the 19 layer VGG-Network (see the architecture in Figure 2). We do not use any of the fully connected layers. The model is publicly available and can be explored in the Caffe-framework. For image synthesis we found that replacing the max-pooling operation by average pooling improves the gradient flow and one obtains slightly more appealing results, which is why the images shown were generated with average pooling.

Generally each layer in the network defines a non-linear filter bank whose complexity increases with the position of the layer in the network. Hence a given input image \vec{x} is encoded in each layer of the CNN by the filter responses to that image. 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 maps. So the responses in a layer l can be stored in a matrix $F^l \in \mathbb{R}^{N_l \times M_l}$ where F_{ij}^l is the activation of the i^{th} filter at position j in layer l . To visualise the image information that is encoded at different layers of the hierarchy we perform gradient descent on a white noise image to find another image that matches the feature responses of the original image. So let \vec{p} and \vec{x} be the original image and the image that is generated and P^l and F^l their respective feature representation in layer l . We then define the squared-error loss between the two feature representations :

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

The derivative of the loss with respect to the activations in layer l equals

$$\frac{\partial L_{content}}{\partial F_{ij}^l} = \begin{cases} (F^l - P^l)_{ij}, F_{ij}^l > 0 \\ 0, otherwise \end{cases} \quad (2)$$

from which the gradient with respect to the image \vec{x} can be computed using standard error back-propagation. Thus we can change the initially random image \vec{x} until it generates the same response in a certain layer of the CNN as the original image \vec{p} .

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 :

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \quad (3)$$

To generate a texture that matches the style of a given image, In the project, we use gradient descent from a white noise image to find another image that matches the style representation of the original image. This is done by minimising the mean-squared distance between the entries of the Gram matrix from the original image and the Gram matrix of the image to be generated. So let \vec{a} and \vec{x} be the original image and the image that is generated and A^l and G^l their respective style representations in layer l . The contribution of that layer to the total loss is then

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

and the total loss is

$$L_{style}(\vec{a}, \vec{x}) = \sum_{l=0}^L \omega_l E_l \quad (5)$$

where ω_l are weighting factors of the contribution of each layer to the total loss. The derivative of E_l with respect to the activations in layer l can be computed by this:

$$\frac{\partial E_l}{\partial F_{ij}^l} = \frac{1}{N_l^2 M_l^2} ((F^l)^T (G^l - A^l))_{ji} \quad (6)$$

The gradients of E_l with respect to the activations in lower layers of the network can be readily computed using standard error back-propagation.

To generate the image that mix the content of a picture with the style of a painting we jointly minimise the distance of a white noise image from the content representation of the photograph in one layer of the network and the

style representation of the painting in a number of the CNN. So let \vec{p} be the picture and \vec{a} be the artwork. The loss function we minimise is :

$$L_{total}(\vec{a}, \vec{p}, \vec{x}) = \alpha L_{content}(\vec{p}, \vec{x}) + \beta L_{style}(\vec{a}, \vec{x}) \quad (7)$$

Where α and β are the weighting factors for content and style reconstruction respectively.

参考文献

- [1] J. Deng, A. Berg, S. Satheesh, H. Su, A. Khosla, and Fei-Fei. Large scale visual recognition challenge. In *In <http://www.image-net.org/challenges/LSVRC/2012/index>*, 2012.
- [2] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *CoRR*, abs/1508.06576, 2015.
- [3] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems*, 25:2012, 2012.
- [5] K. Simonyan, A. Zisserman, K. Simonyan, and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *Eprint Arxiv*, 2014.

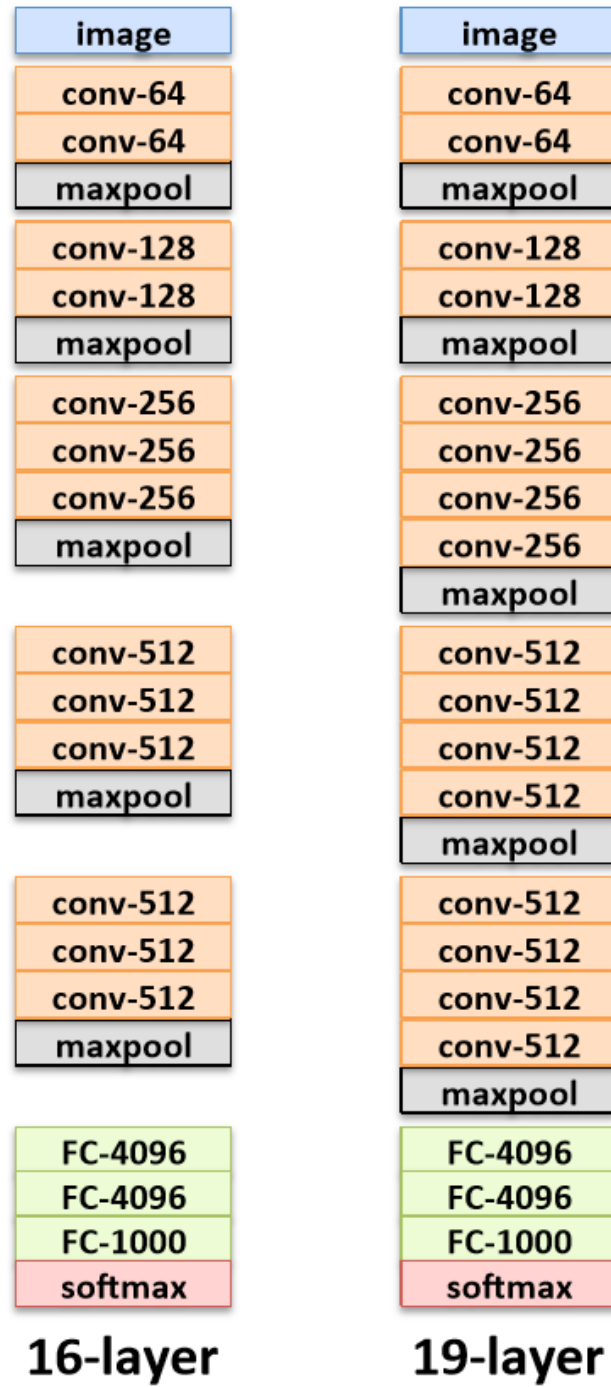


图 2: Graph of a convolutional layer