

# Reading Notes: Higher-order Integration of Hierarchical Convolutional Activations for Fine-grained Visual Categorization

theppsh

November 21, 2018

## 1 Introduction

This paper [1] incorporates the *kernel representations* into Convolutional Neural Networks (CNN) for the fine-grained visual categorization (FGVC). I think this trick can be taken into some object detectors by modifying something... And it might lead to higher performance, so I step into this article for more consideration.

In general, taking the higher-order statistics into CNN by kernel trick can help find the discriminative information between two categories, and benefits the final object classification or detection, and it seems to be fashion that incorporating this traditional method into deep learning recently. Thus, the authors leverage several sets of convolutional filters to get an approximation of  $r$ -order polynomial kernels.

## 2 Matching kernel and polynomial predictor

Assume the output from some specific convolutional layers be tensor  $\mathcal{X} \in \mathbb{R}^{K \times M \times N}$ , and tensor  $\mathcal{X} = \{\mathbf{x}_p\}_{p \in \Omega}$ , where  $\mathbf{x}_p \in \mathbb{R}^K$  represents the descriptor at a particular position  $p$  over the set  $\Omega$  of valid spatial locations with  $|\Omega| = M \times N$ .

### 2.1 Review kernel trick

Kernel tricks are widely used in the conventional machine learning algorithms, for instance, the *kernel PCA*, *kernel fisher analysis*, and *kernel SVM*, etc. The intuitive thought of kernelization is to map the original low-dimensional data to a high-dimensional or infinite-dimensional space by a function  $\phi : \mathbf{x} \in \mathbb{R}^{low} \rightarrow \mathbf{y} \in \mathbb{R}^{high}$ ,  $low \lesssim high$ , thus the non-linear-discriminative data can be converted to be linearly separable. The most fascinating is most methods leverage *kernel trick* to get the pair-wise similarity graph of the original data rather than calculate the high-dimensional data directly, so it's effective and efficient.

Consider the matching scheme  $\mathcal{K}$  for activation sets  $\mathcal{X}$  and  $\bar{\mathcal{X}}$  from two images, in which the set similarity is measured by aggregating all the pairwise similarities among the local descriptors:

$$\mathcal{K}(\mathcal{X}, \bar{\mathcal{X}}) = g(\{k(\mathbf{x}_p, \mathbf{x}_q)\}_{p,q \in \Omega}) = \psi(\mathcal{X})^T \psi(\bar{\mathcal{X}}), \quad (1)$$

where  $k(\cdot)$  is some kernel function to obtain the similarity between  $\mathbf{x}_p$  and  $\mathbf{x}_q$ ,  $g(\cdot)$  is a global sum-pooling function in this paper,  $\psi(\mathcal{X})$  and  $\psi(\bar{\mathcal{X}})$  are the vector representations for the two sets individually. It is worth noting that the construction of  $\psi(\cdot)$  presented above is decomposed into two steps in CNNs: feature mapping and feature aggregating. The mapping step maps each local descriptor  $\mathbf{x} \in \mathbb{R}^K$  into  $\phi(\mathbf{x}) \in \mathbb{R}^D$ . The feature aggregating produces an image-level representations  $\psi(\mathcal{X}) \in \mathbb{R}^D$  by applying global sum-pooling function on the sets  $\{\phi(\mathbf{x}_p)\}_{p \in \Omega}$  after the mapping step.

Note that Eqn.(1) can be rewritten as:

$$\begin{aligned} \mathcal{K}(\mathcal{X}, \bar{\mathcal{X}}) &= \left( \sum_p \phi(\mathbf{x}_p) \right)^T \left( \sum_q \phi(\mathbf{x}_q) \right) \\ &= \sum_p (\phi(\mathbf{x}_p)^T \sum_q \phi(\mathbf{x}_q)) \\ &= \sum_p \sum_q \phi(\mathbf{x}_p)^T \phi(\mathbf{x}_q) \end{aligned} \quad (2)$$

Since we use the global sum-pooling function as  $g(\cdot)$ , thus  $\psi(\mathcal{X})$  satisfies:

$$\psi(\mathcal{X}) = \sum_p \phi(\mathbf{x}_p). \quad (3)$$

The term  $\phi(\mathbf{x}_p)^T \sum_q \phi(\mathbf{x}_q)$  in Eqn.(2) can be interpreted as the linear predictor of  $\phi(\mathbf{x}_p)$ , i.e. we regard  $\sum_q \phi(\mathbf{x}_q)$  as the weights. Inspired by the above linear predictor, we define a function  $f(\cdot)$  to project  $\phi(\mathbf{x}_p) \in \mathbb{R}^D$  to  $\mathbb{R}$ :

$$f(\mathbf{x}_p) = \langle \mathbf{w}, \phi(\mathbf{x}_p) \rangle, \quad (4)$$

where  $\langle \cdot, \cdot \rangle$  is the dot product. To capture the complex and high-order information among parts, use the following polynomial predictor:

$$\begin{aligned} f(\mathbf{x}_p) &= \sum_{k=1}^K w_k x_p^k + \sum_{r=2}^R \sum_{k_1, \dots, k_r} \mathcal{W}_{k_1, \dots, k_r}^r \left( \prod_{s=1}^r x_p^{k_s} \right) \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \langle \mathcal{W}^r, \otimes_r \mathbf{x}_p \rangle \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \sum_{k_1, \dots, k_r} \prod_{s=1}^r \sqrt{\mathcal{W}_{k_1, \dots, k_r}^r} x_p^{k_s} \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \prod_{s=1}^r \langle \sqrt{\mathcal{W}_{k_1, \dots, k_r}^r} \cdot \mathbf{1}, \mathbf{x}_p \rangle \end{aligned} \quad (5)$$

where  $\mathbf{w}$  is the paramter of 1-order predictor,  $\mathcal{W}^r \in \mathbb{R}^{K^r}$  is  $r$ -order tensor which contains the weights of degree- $r$  variable combination in  $\mathbf{x}$ ,  $\otimes_r \mathbf{x} = \underbrace{\mathbf{x} \otimes \dots \otimes \mathbf{x}}_r$ ,

and  $\otimes$  denotes the *kronecker product*.  $\mathbf{1} \in \mathbb{R}^K$  is a vector whose entries equal to 1. Note that function  $f(\cdot)$  is just a mapping for  $\mathbf{x}$ , it shares some ideas like kernel function in Eqn.(1), but actually they are totally different. Obviously, take an example of the equation:

$$\sum_p \sum_q w_p x_p w_q x_q = \langle \mathbf{w}, \mathbf{x} \rangle^2 = \prod_{i=1}^2 \langle \mathbf{w}, \mathbf{x} \rangle$$

## 2.2 Tensor learning

The  $r$ -order tensor  $\mathbf{W}^r \in \mathbb{R}^{K^r}$  in Eqn.(5) can be approximated by rank-one tensor decomposition: Assume  $r$  vectors  $\mathbf{u}_1, \dots, \mathbf{u}_{K_1} \in \mathbb{R}^{K_r}$ , and their outer product is  $K_1 \times \dots \times K_r$  rank-one tensor which satisfies  $(\mathbf{u}_1 \otimes \dots \otimes \mathbf{u}_r)_{k_1, \dots, k_r} = (\mathbf{u}_1)_{k_1} \dots (\mathbf{u}_r)_{k_r}$ . The rank-one decomposition for such tensor  $\mathbf{W}$  is defined as  $\mathbf{W} = \sum_{d=1}^D \alpha^d \mathbf{u}_1^d \otimes \dots \otimes \mathbf{u}_r^d$ , where  $\alpha^d$  is the weight for  $d$ -th rank-one tensor,  $D$  is the rank of tensor. We then apply the rank-one approximation for each  $r$ -order tensor  $\mathbf{W}^r$  and present the following alternative form of polynomial predictor in Eqn.(5):

$$\begin{aligned} f(\mathbf{x}_p) &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \left\langle \sum_{d=1}^{D^r} \alpha^{r,d} \mathbf{u}_1^{r,d} \otimes \dots \otimes \mathbf{u}_r^{r,d}, \otimes_r \mathbf{x}_p \right\rangle \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \sum_{d=1}^{D^r} \alpha^{r,d} \langle \mathbf{u}_1^{r,d} \otimes \dots \otimes \mathbf{u}_r^{r,d}, \otimes_r \mathbf{x}_p \rangle \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \sum_{d=1}^{D^r} \alpha^{r,d} \sum_{k_1, \dots, k_r} \prod_s^r u_{s,k_s}^{r,d} x_p^{k_s} \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \sum_{d=1}^{D^r} \alpha^{r,d} \prod_s^r \langle \mathbf{u}_s^{r,d}, \mathbf{x}_p \rangle \\ &= \langle \mathbf{w}, \mathbf{x}_p \rangle + \sum_{r=2}^R \langle \boldsymbol{\alpha}^r, \mathbf{z}_p^r \rangle \end{aligned} \quad (6)$$

where  $\mathbf{u}_1^{r,d}, \dots, \mathbf{u}_r^{r,d} \in \mathbb{R}^K$  are parameters can be learned end-to-end.  $D^r$  represents the minimal rank of the  $r$ -order tensor  $\mathbf{W}^r$ , or the hyper-parameter corresponding to the output channels of a particular CNN.  $\boldsymbol{\alpha}^r = [\alpha^{r,1}, \dots, \alpha^{r,D^r}]^T$  is the associated weight vector of all  $D^r$  rank-one tensors. Vector  $\mathbf{z}_p^r \in \mathbb{R}^{D^r}$  whose  $d$ -th element satisfies:

$$(\mathbf{z}_p^r)_d = \prod_{s=1}^r \langle \mathbf{u}_s^{r,d}, \mathbf{x}_p \rangle \quad (7)$$

which characterizes the degree- $r$  variable interactions under a single rank-one tensor basis.

That paper costs a lot of time to describe the final function  $f(\mathbf{x}_p)$  in Eqn.(6), and I think they will implement  $f(\cdot)$  in CNN fashion, but it seems they just use  $\mathbf{z}_p^r$  incorporated in the final models.

The implementation embedded in CNN can be described as follows: firstly, for each  $s$ , we can deploy  $\{\mathbf{u}_s^{r,d}\}_{d=1, \dots, D^r}$  as a set of  $D^r$   $1 \times 1$  convolutional

filters to generate the set of feature maps  $\mathcal{Z}_s^r$  of dimension  $D^r \times M \times N$ . To get the high-order information, we apply element-wise multiplication denoted as  $\odot$ , and get the output  $\mathcal{Z}^r \in \mathbb{R}^{D^r \times M \times N}$  with high-order statistics information as follows:

$$\mathcal{Z}^r = \mathcal{Z}_1^r \odot \cdots \odot \mathcal{Z}_r^r, \quad (8)$$

then sum-pooling function  $g(\cdot)$  is used to generate the final output  $\mathbf{y}^r \in \mathbb{R}^{D^r}$  by sum-up all the local descriptors in  $\mathcal{Z}^r$ :

$$\mathbf{y}^r = g(\mathcal{Z}^r) = \sum_{p \in \Omega} \mathbf{z}_p^r, \quad (9)$$

referring to Eqn.(9), the derivatives for  $\mathbf{x}_p$  and each degree- $r$  convolutional filter  $\mathbf{u}_s^{r,d}$  in back propagation process can be achieved by:

$$\frac{\partial \ell}{\partial \mathbf{z}_p^r} = \frac{\partial \ell}{\partial \mathbf{y}^r} \quad (10)$$

from Eqn.(7) and Eqn.(8), we denote  $\mathbf{z}_p^r, p \in \Omega$  as:

$$\mathbf{z}_p^r = \mathbf{z}_{p,1}^r \odot \cdots \odot \mathbf{z}_{p,r}^r, \quad (11)$$

where we define  $\mathbf{z}_{p,s}^r \in \mathbb{R}^{D^r}$  whose  $d$ -th entry is:

$$(\mathbf{z}_{p,s}^r)_d = \langle \mathbf{u}_s^{r,d}, \mathbf{x}_p \rangle, \quad (12)$$

notice the symbol  $\mathbf{z}_{p,d}^r$  is different from term  $(\mathbf{z}_p^r)_d$  in Eqn.(7). Thus the derivative of loss  $\ell$  respect to  $\mathbf{z}_{p,s}^r$  is:

$$\frac{\partial \ell}{\partial \mathbf{z}_{p,s}^r} = \frac{\partial \ell}{\partial \mathbf{z}_p^r} \odot \frac{\partial \mathbf{z}_p^r}{\partial \mathbf{z}_{p,s}^r} = \frac{\partial \ell}{\partial \mathbf{z}_p^r} \odot \mathbf{z}_{p,1}^r \odot \cdots \odot \mathbf{z}_{p,s-1}^r \odot \mathbf{z}_{p,s+1}^r \odot \cdots \odot \mathbf{z}_{p,r}^r \quad (13)$$

For the later convenience of derivation, define a matrix  $\mathbf{U}_s^r = [\mathbf{u}_s^{r,1}, \dots, \mathbf{u}_s^{r,D^r}]^T \in \mathbb{R}^{D^r \times K}$ . Thus

$$\mathbf{z}_{p,s}^r = \mathbf{U}_s^r \mathbf{x}_p, \quad (14)$$

therefore, the derivative of  $\ell$  respect to  $\mathbf{x}_p$  can be obtained by:

$$\frac{\partial \ell}{\partial \mathbf{x}_p} = \frac{\partial \ell}{\partial \mathbf{y}^1} + \sum_{r=2}^R \sum_{s=1}^r (\mathbf{U}_s^r)^T \frac{\partial \ell}{\partial \mathbf{z}_{p,s}^r}, \quad (15)$$

notice that we regard  $\mathbf{z}_p^1$  as  $\mathbf{x}_p$ , thus the term  $\frac{\partial \ell}{\partial \mathbf{y}^1}$  in Eqn.(15) represents the derivative respect to order-one interactions in  $\mathbf{x}_p$ . Use the chain rules, the same procedure can be easily adopted into the derivative respect to  $\mathbf{U}_s^r$ :

$$\frac{\partial \ell}{\partial \mathbf{U}_s^r} = \sum_{p \in \Omega} \frac{\partial \ell}{\partial \mathbf{z}_{p,s}^r} (\mathbf{x}_p)^T \quad (16)$$

Of course, all the operation can be easily extended to matrix multiplication rather than single sample  $\mathbf{x}_p$ , you just need to be careful when programing. Actually, It is worth noting that some frameworks such PyTorch have implemented the derivative procedure.

### 3 Experiments

To be continued...

### References

- [1] Sijia Cai, Wangmeng Zuo, and Lei Zhang. Higher-order integration of hierarchical convolutional activations for fine-grained visual categorization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 511–520, 2017.