

Learning a Robust Representation via a Deep Network on Symmetric Positive Definite Manifolds

Zhi Gao, Yuwei Wu, Xingyuan Bu, and Yunde Jia *Member, IEEE*

Abstract—Recent studies have shown that aggregating convolutional features of a pre-trained Convolutional Neural Network (CNN) can obtain impressive performance for a variety of visual tasks. The symmetric Positive Definite (SPD) matrix becomes a powerful tool due to its remarkable ability to learn an appropriate statistic representation to characterize the underlying structure of visual features. In this paper, we propose to aggregate deep convolutional features into an SPD matrix representation through the SPD generation and the SPD transformation under an end-to-end deep network. To this end, several new layers are introduced in our network, including a nonlinear kernel aggregation layer, an SPD matrix transformation layer, and a vectorization layer. The nonlinear kernel aggregation layer is employed to aggregate the convolutional features into a real SPD matrix directly. The SPD matrix transformation layer is designed to construct a more compact and discriminative SPD representation. The vectorization and normalization operations are performed in the vectorization layer for reducing the redundancy and accelerating the convergence. The SPD matrix in our network can be considered as a mid-level representation bridging convolutional features and high-level semantic features. To demonstrate the effectiveness of our method, we conduct extensive experiments on visual classification. Experiment results show that our method notably outperforms state-of-the-art methods.

Index Terms—Feature Aggregation, SPD Matrix, Riemannian Manifold, Deep Convolutional Network

I. INTRODUCTION

DEEP Convolutional Neural Networks (CNNs) have shown great success in many vision tasks. There are several successful networks, *e.g.*, AlexNet [1], VGG [2], GoogleNet [3], Network In Network [4] and ResNet [5]. Driven by the emergence of large-scale data sets and fast development of computation power, features based on CNNs have proven to perform remarkably well on a wide range of visual recognition tasks [6], [7]. Two contemporaneous works introduced by Liu *et al.* [8] and Babenko and Lempitsky [9] demonstrate that convolutional features could be seen as a set of local features which can capture the visual representation related to objects. To make better use of deep convolutional features, many efforts have been devoted to aggregating them, such as max pooling [10], cross-dimensional pooling [11], sum pooling [9], and bilinear pooling [8], [12]. However, modeling these convolutional features to boost the feature learning ability of a CNN is still a challenging task. This work investigates a more effective scheme to aggregate convolutional features

to produce a robust representation using an end-to-end deep network for visual tasks.

Recently, the Symmetric Positive Definite (SPD) matrix has been demonstrated the powerful representation ability and widely used in computer vision community, such as face recognition [13], [14], image set classification [15], transfer learning [16], and action recognition [17], [18]. Through the theory of non-Euclidean Riemannian geometry, the SPD matrix often turns out to be better suited in capturing desirable data distribution properties. Accordingly, we attempt to aggregate the deep convolutional features into a powerful SPD matrix as a robust representation.

The second-order statistic information of convolutional features, *e.g.*, the covariance matrix and Gaussian distribution, are the widely used SPD matrix representation endowed with CNNs [19], [20], [21]. The dimensionality of convolutional features extracted from CNNs may be much larger than that of hand-craft features. As a result, modeling convolutional features from CNNs by using the covariance matrix or Gaussian distribution is insufficient to precisely model the real feature distribution. When the dimension of features is larger than the number of features, the covariance matrix and Gaussian distribution is a symmetric Positive SemiDefinite (PSD) matrix, *i.e.*, the singular matrix. Singular matrix makes the data have an unreasonable manifold structure. In this case, the Riemannian metrics, *e.g.*, the affine-invariant distance and Log-Euclidean distance, are unsuitable to measure the manifold structure of SPD matrices. Moreover, most SPD matrix embedding on deep networks only contains the linear correlation of features. Owing the ability of capturing nonlinear relationship among features is indispensable for a generic representation.

It is thus desirable that a more discriminative and suitable SPD representation aggregated from deep features should be established in an end-to-end framework for visual analysis. To this end, we design a series of new layers to overcome existing issues aforementioned based on the following two observations.

- Kernel functions possess an ability of modeling nonlinear relationships of data, and they are flexible and easy to be computed. Beyond covariance [22] have witnessed significant advances of positive definite kernel functions whose kernel matrices are real SPD matrices, no matter what the number of the feature dimension and the number of features are. Since many kernel functions are differentiable, such as Radial Basis Function (RBF) kernel function, Polynomial kernel function and Laplacian kernel function [23], they can be readily embedded into a network to implement an end-to-end training, which

The authors are with Beijing Laboratory of Intelligent Information Technology, School of Computer Science, Beijing Institute of Technology (BIT), Beijing, China. Email: {gaozhi_2017, wuyuwei, buxingyuan, jiayunde}@bit.edu.cn. Corresponding author: Yuwei Wu.

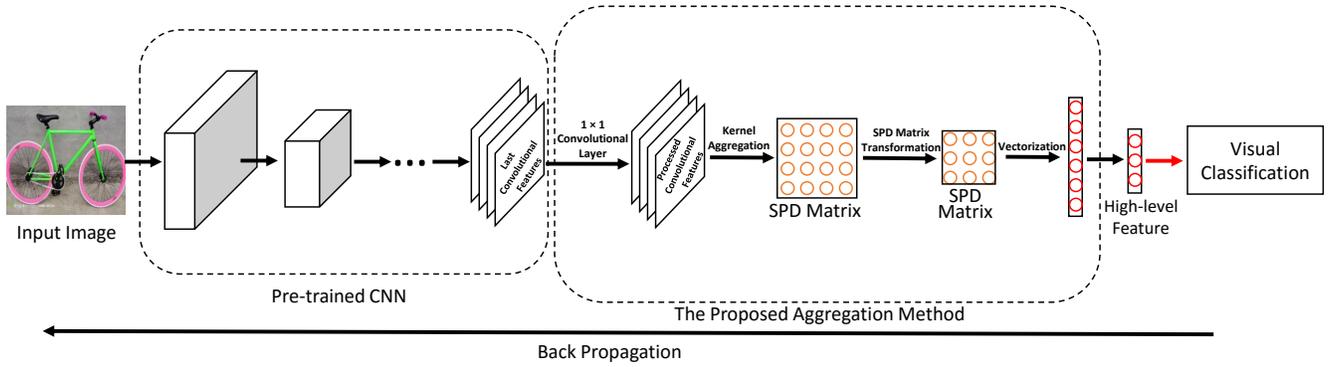


Fig. 1. The flowchart of our SPD aggregation network. We focus on both how to generate an SPD matrix from convolutional features and gain a more discriminative SPD representation by the transformation operation. Our general representation can be used for different tasks. To gain the processed convolutional features, the PCA and 1×1 convolutional layer are applied on the last convolutional features from the pre-trained CNN network. The kernel aggregation operation, SPD matrix transformation operation and vectorization operation correspond to the proposed three new layers. The kernel aggregation operation is to generate an SPD representation from the convolutional features. The SPD matrix transformation operation is to map the SPD matrix to a more compact and discriminative one by learnable parameters. The vectorization operation involves vectorization and normalization operations on the transformed SPD matrix.

is well aligned with the design requirements of a deep network.

- Several deep SPD networks [24], [25], [26] transform the SPD matrix to a new compact and discriminative matrix. The network input is an SPD matrix. The transformed matrix is still an SPD matrix which can capture desirable data properties. We find that the transformed SPD matrix after the learnable layers leads to better performance than the original SPD matrix. The output SPD matrix not only have characteristics of a general SPD matrix that captures the desirable properties of visual features but also is more suitable and discriminative to the specific visual task.

Motivated by empirical observations mentioned above, we introduce a convolutional feature aggregation operation which consists of the SPD generation and the SPD transformation. Three new layers including a kernel aggregation layer, an SPD matrix transformation layer and a vectorization layer, are designed to replace the traditional global pooling layers and fully connected (FC) layers. Concretely, we deem each feature map as a sample and present a kernel aggregation layer using a nonlinear kernel function to generate an SPD matrix. The proposed kernel matrix models a nonlinear relationship among feature maps and ensures that the SPD matrix is nonsingular. More importantly, our kernel matrix is differentiable, which entirely meets requirements of a deep network. The SPD matrix transformation layer is employed to map the SPD matrix to a more discriminative and compact one. Thanks to the symmetry property, the vectorization layer carries out the upper triangle vectorization and normalization operations to the SPD matrix followed by the classifier. The architecture of our network is illustrated in Fig. 1. The proposed method first generates an SPD matrix based on convolutional features and then transforms the initial SPD matrix to a more discriminative one. It can not only capture the real spatial information but also encode high-level variation information among convolutional features. Actually, the obtained descriptor acts as a mid-level representation bridging convolutional features and high-level semantics features. The resulting vector can contribute to visual classification tasks, as validated in experiments.

In summary, our contributions are three-fold.

- (1) We apply the SPD matrix non-linear aggregation to the convolutional feature aggregation field by the generation and the transformation two processes. In this way, it can learn a compactness and robustness SPD matrix representation to characterize the underlying structure of convolutional features.
- (2) We carry out the nonlinear aggregation of convolutional features under a Riemannian deep network architecture, where three novel layers are introduced, *i.e.*, a kernel aggregation layer, an SPD matrix transformation layer and a vectorization layer. The state-of-the-art performance of our SPD aggregation network is consistently achieved over the visual classification tasks.
- (3) We exploit the faster matrix operations to avoid the cyclic calculation in forward and backward backpropagations of the kernel aggregation layer. In addition, we present the component decomposition and retraction of the Orthogonal Stiefle manifold to carry out the backpropagation on the SPD matrix transformation layer.

The remaining sections are organized as follows. We review the recent works about feature aggregation methods in both Euclidean Space and Riemannian Space in Section II. Section III presents the details of our SPD aggregation method. We report and discuss the experimental results in Section IV, and conclude the paper in Section V.

II. RELATED WORK

Feature aggregation is an important problem in computer vision tasks. Recent works have witnessed significant advances of CNNs. It is still a challenging work to find a suitable way to aggregate convolutional features. In this section, we review typical techniques of feature aggregation in both the Euclidean space and Riemannian space.

A. Convolutional Feature Aggregation in the Euclidean Space

An effective image representation is an essential element for visual recognition due to the object appearance variations caused by pose, view angle, and illumination changes.

Traditional methods typically obtain the image representation by aggregating hand-crafted local features (*e.g.*, SIFT) into a global image descriptor. Popular aggregation schemes include Bag-of-words (BOW) [27], Fisher Vector (FV) [28], and Vector of Locally Aggregated Descriptor (VLAD) [29]. Gong *et al.* [30] introduced a multi-scale orderless pooling scheme to aggregates FC6 features of local patches into a global feature using VLAD. The VLAD ignores different effects of each cluster center. Cimpoi *et al.* [31] treated the convolutional layer of CNNs as a filter bank and built an orderless representation using FV. In addition, Liu *et al.* [32] proposed a cross convolutional layer pooling scheme which regards feature maps as a weighting filter to the local features. Toliás *et al.* [10] max pooled convolutional features of the last convolutional layer to represent each patch and achieved compelling performance for object retrieval. Babenko *et al.* [9] compared different kinds of aggregation methods (*i.e.*, max pooling, sum pooling and fisher vector) for last convolutional layer features and demonstrated the sum-pooled convolutional descriptor is really competitive with other aggregation schemes.

Works mentioned above only treat the CNN as a black-box feature extractor rather than studying on properties of CNN features in an end-to-end framework. Several researchers [8], [33], [34] suggested that the end-to-end network can achieve better performance because it is sufficient by itself to discover good features for visual tasks. Arandjelovic *et al.* [34] proposed a NetVLAD which adopts an the end-to-end framework for weakly supervised place recognition. Based on the ResNet, Zhang *et al.* [33] introduced an extended version of the VLAD, *i.e.*, Deep-TEN, for texture classification. Lin *et al.* [8] presented a general orderless pooling model named Bilinear to compute the outer product of local features. He *et al.* [35] introduced a spatial pyramid pooling method eliminating the constrain of the fixed-size input image.

Recent research shows that exploiting the manifold structure representation is more effective than the hypothetical Euclidean distribution in some visual tasks. The difference between our method and the traditional aggregation methods in the Euclidean space is that we use the powerful SPD manifold structure to aggregate the desirable data distributions of features. We design an SPD aggregation scheme to generate the SPD matrix as the resulting representation, and transform the SPD representation to more discriminative one by learnable layers.

B. Convolutional Feature Aggregation in the Riemannian Space

The aggregation methods in the non-Euclidean space have been successful applied. It can capture more appropriate feature distributions information. The second-order statistic information has better performance than the first-order statistic [19], such as average pooling. Some works directly regard the second-order statistic information as the SPD matrix. Ionescu *et al.* [20] proposed a DeepO2P network that uses a covariance matrix as the image representation. They mapped points on the manifold to the logarithm tangent space and derived a new chain rule for derivatives. Li *et al.* [19] presented a

matrix normalized covariance method exploring the second-order statistic. This work can tackle the singular issue of the covariance matrix by the normalization operation. Yu and Salzmann [21] introduced a covariance descriptor unit to integrates second-order statistic information. The covariance matrix of convolutional features is generated and then transformed to a vector for the softmax classifier. Compared with our network, these three works are confined to the drawbacks of covariance matrices. Engin *et al.* [36] designed a deep kernel matrix based SPD representation, but didn't contains the transformation process.

Other SPD Riemannian networks mainly project an SPD matrix to a more discriminative one. Dong *et al.* [24] and Huang and Gool [25] proposed Riemannian networks contemporaneously, in which the inputs of their networks are SPD matrices. The networks projects high dimensional SPD matrices to a low dimensional discriminative SPD manifold by a nonlinear mapping. Zhang *et al.* [26] introduced new layers to transform and vectorize the SPD matrix for action recognition, where the input is a nonlinear kernel matrix modeling correlation of frames in a video. However, these three works only focused on how to transform the SPD matrix without utilizing the powerful convolutional features. The generation of the input SPD matrix can not be guided by the loss function. In contrast, our method focuses on not only the SPD matrix transformation but also the generation from convolutional features.

Our work is closely related with [19], [21], [36]. We make it clear that the proposed convolutional feature aggregation method is composed of generation and transformation processes. Compared with [19], our method utilizes the kernel matrix as the representation instead of the second-order statistic covariance matrix, characterizing complex nonlinear variation information of features. In addition, our aggregation method contains a learnable transformation process than [19], making SPD representation more compact and robust. The generated SPD matrix in our method is more powerful than the covariance matrix in [21], avoiding some drawbacks of PSD matrix. In addition, instead of a transformation from a matrix to a vector, the vectorization operation in our work is taking the upper triangle of a matrix since there are already transformation operations between the SPD matrices. Compared to [36], our SPD representation can be more compact and robust through the transformation process.

III. SPD AGGREGATION METHOD

Our model aims to aggregate convolutional features into a powerful SPD representation in an end-to-end fashion. To this end, we design three novel layers including a kernel aggregation layer, an SPD matrix transformation layer and a vectorization layer. Our SPD aggregation can be applied to the visual classification. Specifically, the convolutional features pass through the proposed three layers followed by an FC layer and a loss function. The intermediate generated SPD matrix can be treated as a mid-level representation which is a connection between convolutional features and high-level features. The architecture of our network is illustrated in Fig. 1(c).

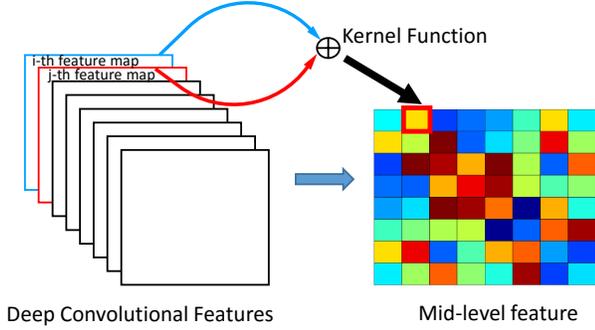


Fig. 2. Using the kernel function to generate an SPD matrix from convolutional features. The input is convolutional features and the output is an SPD kernel matrix.

A. Preprocessing of Convolutional Features

A CNN model trained on a large dataset such as ImageNet can have a better general representation ability. We would like to fuse the convolutional features of the last convolutional layer and adjust the dimension of convolutional features for different tasks. We introduce a convolutional layer whose filter's size is 1×1 between the last convolutional layer of the off-the-shelf model and the kernel aggregation layer to make the processed convolutional features more adaptive to the SPD matrix representation. A Relu layer follows the 1×1 convolutional layer to increase the nonlinear ability.

B. Kernel Aggregation Layer

We present the kernel aggregation layer to aggregate convolutional features into an initial SPD matrix. Let $X \in \mathbb{R}^{C \times H \times W}$ be 3-dimensional convolutional features. C is the number of channels, *i.e.*, the number of feature maps, H and W are the height and width of each feature map, respectively. Let $x_i \in \mathbb{R}^C$ denote the i -th local feature, and there are N local features in total, where $N = H \times W$. $f_i \in \mathbb{R}^{H \times W}$ is the i -th feature map.

Although several approaches have applied a covariance matrix Cov to be a generic feature representation and obtained promising results, two issues remain to be addressed. First, the rank of covariance matrix should hold $rank(Cov) \leq \min(C, N - 1)$, otherwise covariance matrix is prone to be singular when the dimension C of local features is larger than the number of local features extracted from an image region. Second, for a generic representation, the capability of modeling nonlinear feature relationship is essential. However, covariance matrix only evaluates the linear correlation of features.

To address these issues, we adopt the nonlinear kernel matrix as a generic feature representation to aggregate deep convolutional features. In particular, we take advantage of the Riemannian structure of SPD matrices to describe the second-order statistic and nonlinear correlations among deep convolutional features. The nonlinear kernel matrix is capable of modeling nonlinear feature relationship and is guaranteed

to be nonsingular. Different from the traditional kernel-based methods whose entries evaluates the similarity between a pair of samples, we apply the kernel mapping to each feature f_1, f_2, \dots, f_C rather than each sample x_1, x_2, \dots, x_N . Mercer kernels are usually employed to carry out the mapping implicitly. The Mercer kernel is a function $\mathcal{K}(\cdot, \cdot)$ which can generate a kernel matrix $K \in \mathbb{R}^{C \times C}$ using pairwise inner products between mapped convolutional features for all the input data points. The K_{ij} in our nonlinear kernel matrix K can be defined as

$$K_{ij} = \mathcal{K}(f_i, f_j) = \langle \phi(f_i), \phi(f_j) \rangle, \quad (1)$$

where $\phi(\cdot)$ is an implicit mapping. In this paper, we exploit the Radial Basis Function (RBF) kernel function expressed as

$$\mathcal{K}(f_i, f_j) = \exp(-\|f_i - f_j\|^2 / 2\sigma^2), \quad (2)$$

where σ is a positive constant and set to the mean Euclidean distances of all feature maps. What Eq. (2) reveals is the nonlinear relationship between convolutional features.

We show an important theorem for kernel aggregation operation. Based on the Theorem 1, the kernel matrix K of the RBF kernel function is guaranteed to be positive definite no matter what C and N are.

Theorem 1. Let $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^M$ denotes a set of different points and $\mathbf{x}_i \in \mathbb{R}^n$. Then the kernel matrix $K \in \mathbb{R}^{M \times M}$ of the RBF kernel function \mathcal{K} on \mathbf{X} is guaranteed to be a positive definite matrix, whose (j, k) -th element is $K_{jk} = \mathcal{K}(\mathbf{x}_j, \mathbf{x}_k) = \exp(-\alpha\|\mathbf{x}_j - \mathbf{x}_k\|^2)$ and $\alpha > 0$.

Proof. The Fourier transform convention $\hat{\mathcal{K}}(\xi)$ of the RBF kernel function $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\alpha\|\mathbf{x}_j - \mathbf{x}_k\|^2)$ is

$$\hat{\mathcal{K}}(\xi) = (2\pi/\alpha)^{n/2} \int_{\mathbb{R}^n} e^{i\xi\mathbf{x}_j} e^{-i\xi\mathbf{x}_k} e^{-\|\xi\|^2/(2\alpha)} d\xi. \quad (3)$$

Then we calculate the quadratic form of the kernel matrix K . Let $\mathbf{c} = (c_1, \dots, c_M) \in \mathbb{R}^{M \times 1}$ denote an arbitrary nonzero vector. The quadratic form Q is

$$\begin{aligned} Q &= \mathbf{c}^\top K \mathbf{c} = \sum_{j=1}^M \sum_{k=1}^M c_j c_k \exp(-\alpha\|\mathbf{x}_j - \mathbf{x}_k\|^2) \\ &= \sum_{j=1}^M \sum_{k=1}^M c_j c_k (2\pi/\alpha)^{n/2} \int_{\mathbb{R}^n} e^{i\xi\mathbf{x}_j} e^{-i\xi\mathbf{x}_k} e^{-\|\xi\|^2/(2\alpha)} d\xi \\ &= (2\pi/\alpha)^{n/2} \int_{\mathbb{R}^n} e^{-\|\xi\|^2/(2\alpha)} \left\| \sum_{k=1}^M c_k e^{-i\xi\mathbf{x}_k} \right\|^2 d\xi, \end{aligned} \quad (4)$$

where \top is the transpose operation. Because $e^{-\|\xi\|^2/(2\alpha)}$ is a positive and continuous function, the quadratic form $Q = 0$ on the condition that

$$\sum_{k=1}^M c_k e^{-i\xi\mathbf{x}_k} = 0. \quad (5)$$

However, the complex exponentials $e^{-i\xi\mathbf{x}_1}, \dots, e^{-i\xi\mathbf{x}_M}$ is linear independence. Accordingly, $Q > 0$ and kernel matrix K is a positive definite matrix. \square

In this work, K is the generated SPD matrix as the mid-level image representation. Any SPD manifold optimization can be applied directly, without structure being destroyed. The toy example of the kernel aggregation is illustrated in Fig. 2. As we all known, the kernel aggregation layer should be differentiable to meet the requirement of an end-to-end deep learning framework. Clearly, Eq. (2) is differentiable with respect to the input X . Denoting by \mathcal{L} the loss function, the gradient with respect to the kernel matrix is $\frac{\partial \mathcal{L}}{\partial K} \cdot \frac{\partial \mathcal{L}}{\partial K_{ij}}$ is an element in $\frac{\partial \mathcal{L}}{\partial K}$. We compute the partial derivatives of \mathcal{L} with respect to f_i and f_j , which are

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial f_i} &= \sum_{j=1}^{H \times W} \frac{\partial \mathcal{L}}{\partial K_{ij}} \frac{K_{ij}}{-2\sigma^2} (f_i - f_j) \\ \frac{\partial \mathcal{L}}{\partial f_j} &= \sum_{i=1}^{H \times W} \frac{\partial \mathcal{L}}{\partial K_{ij}} \frac{K_{ij}}{-2\sigma^2} (f_j - f_i). \end{aligned} \quad (6)$$

In this process, the gradient of the SPD matrix can flow back to convolutional features.

During forward propagation Eq. (2) and backward propagation Eq. (6), we have to do C^2 cycles to compute the kernel matrix K and $2C^2$ cycles to gain the gradient with respect to convolutional features $\frac{\partial \mathcal{L}}{\partial X}$, where C is the number of channels. Obviously, both the forward and backward propagations are computationally demanding. It is well known that the computation using matrix operations is preferable due to the parallel computing in computers. Accordingly, our kernel aggregation layer is able to be calculated in a faster way via matrix operations. Let's reshape the convolutional features $X \in \mathbb{R}^{C \times H \times W}$ to a matrix $M \in \mathbb{R}^{C \times N}$. Each row of M is a reshaped feature map $\mathbf{f}_i \in \mathbb{R}^{1 \times N}$ obtained from f_i and each column of M is the convolutional local feature x_i . Note that, $\|f_i - f_j\|^2$ in Eq. (2) can be expanded to $\|f_i - f_j\|^2 = \mathbf{f}_i \mathbf{f}_i^\top - 2\mathbf{f}_i \mathbf{f}_j^\top + \mathbf{f}_j \mathbf{f}_j^\top$. For each of inner products $\mathbf{f}_i \mathbf{f}_i^\top$, $2\mathbf{f}_i \mathbf{f}_j^\top$ and $\mathbf{f}_j \mathbf{f}_j^\top$, it needs to be calculated C^2 times in cycles of Eq. (2). Now, we can convert C^2 times inner products operation to a matrix multiplication operation which only needs to be computed once,

$$\begin{aligned} K1 &= \mathbf{1} (M \circ M)^\top \\ K2 &= (M \circ M) \mathbf{1}^\top \\ K3 &= M M^\top, \end{aligned} \quad (7)$$

where \circ is the Hadamard product and $\mathbf{1} \in \mathbb{R}^{C \times N}$ is a matrix whose elements are all "1"s. $K1$, $K2$ and $K3$ are all $C \times C$ real matrices. The element $K1(i, j)$ is the 2-norm of i -th row vector of M , and is equal to the calculation output of $\mathbf{f}_i \mathbf{f}_i^\top$. The element $K2(i, j)$ is the 2-norm of j -th column vector of M , and is equal to the calculation output of $\mathbf{f}_j \mathbf{f}_j^\top$. The element $K3(i, j)$ is equal to $\mathbf{f}_i \mathbf{f}_j^\top$. $K1$, $K2$ and $K3$ can be calculated in advance.

Therefore, we compute $-\|f_i - f_j\|^2/2\sigma^2$ in Eq. (2) by the matrix addition and multiplication, and implement the $\exp(\cdot)$ to the matrix in a parallel computing way instead of calculating each element in the cycle. Then the kernel matrix K can be calculated by matrix operations as follows.

$$K = \exp\left(- (K1 + K2 - 2K3) / 2\sigma^2\right), \quad (8)$$

where $\exp(A)$ means the exponential operation to each element in the matrix A . Although calculating directly the $\exp(\cdot)$ function is time-consuming, it can be computed efficiently in a matrix form through Eq. (8), which is faster than through Eq. (2). Similarly, back propagation process in Eq. (6) can also be carried out in the matrix operation which is given by

$$\frac{\partial \mathcal{L}}{\partial M} = 4 \left(\mathbf{1}^\top \left(\frac{\partial \mathcal{L}}{\partial K} \circ K \right) \circ M^\top - M^\top \left(\frac{\partial \mathcal{L}}{\partial K} \circ K \right) \right)^\top. \quad (9)$$

Remark: The covariance matrix descriptor, as a special case of SPD matrices, captures feature correlations compactly in an object region, and therefore has been proven to be effective for many applications. Given the local features x_1, x_2, \dots, x_N , the covariance descriptor Cov is defined as

$$Cov = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^\top, \quad (10)$$

where $\mu = \frac{1}{N} \sum_{i=1}^N x_i$ is the mean vector. The covariance matrix can also be seen as a kernel matrix where the (i, j) -th element of the covariance matrix Cov can be expressed as

$$Cov_{ij} = \left\langle \frac{\bar{f}_i}{\sqrt{N-1}}, \frac{\bar{f}_j}{\sqrt{N-1}} \right\rangle, \quad (11)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product, $\bar{f}_i = f_i - \mu_i \mathbf{1}$ and μ_i is the mean value of f_i . Therefore, the covariance matrix corresponds to a special case of the nonlinear kernel matrix defined in Eq. (1), where $\phi(f_i) = (\bar{f}_i - \mu_i \mathbf{1}) / \sqrt{N-1}$. Through this way, we can find that covariance matrices contain the simple linear correlation features. Whether the covariance matrix is a positive definite matrix depends on the C and N , i.e., $\text{rank}(Cov) \leq \min(C, N-1)$.

C. SPD Matrix Transformation Layer

As discussed in [24], [25], [26], SPD matrix transformation networks are capable of achieving the better performance than the original SPD matrix. Inspired by [37] and [21], we add a learnable layer to make the network more flexible and more adaptive to the specific task. Based on the SPD matrix generated by the kernel aggregation layer, we expect to transform the existing SPD representation to be a more discriminative, suitable and desirable matrix. To preserve the powerful ability of the SPD matrix, the transformed matrix should also be an SPD matrix. Moreover, we attempt to adjust the dimension to make the SPD matrix more flexible and compact. Here, we design the SPD matrix transformation layer in our network.

Let's define the Riemannian manifold of $n \times n$ SPD matrices as Sym_n^+ . The output SPD matrix K of the kernel aggregation layer lies on the manifold Sym_C^+ . We use a matrix mapping to complete the transformation operation. As depicted in Fig. 3, we map the input SPD matrix which lies on the original manifold Sym_C^+ to a new discriminative and compact SPD

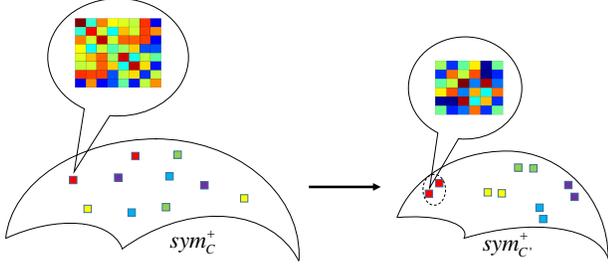


Fig. 3. The illustration of the projection from a manifold to another.

matrix in another manifold $Sym_{C'}^+$, where C' is the dimension of the SPD matrix transformation layer. In this way, the desired transformed matrix can be obtained by a learnable mapping. Given a $C \times C$ SPD matrix K as an input, the output SPD matrix can be calculated as

$$Y = W^\top KW, \quad (12)$$

where $Y \in \mathbb{R}^{C' \times C'}$ is the output of the transformation layer, and $W \in \mathbb{R}^{C \times C'}$ are learnable parameters which are randomly initialized during training. C' controls the size of Y . Based on the Theorem 2, the learnable parameters W should be a column full rank matrix to make Y be an SPD matrix as well.

Theorem 2. Let $A \in \mathbb{R}^{C \times C}$ denote an SPD matrix, $W \in \mathbb{R}^{C \times C'}$ and $B = W^\top AW$, where $C \geq C'$. B is an SPD matrix if and only if W is a column full rank matrix, i.e., $\text{rank}(W) = C'$.

Proof. If A is an SPD matrix, W is a column full rank matrix and $\text{rank}(W) = C'$. For homogeneous equations $W\mathbf{x} = \mathbf{0}$ and $\mathbf{x} \in \mathbb{R}^{C' \times 1}$, $W\mathbf{x} = \mathbf{0}$ only has a zero solution, where $\mathbf{0}$ is the zero vector. For arbitrary nonzero vector \mathbf{x} , $W\mathbf{x} \neq \mathbf{0}$. We calculate the quadric form $\mathbf{x}^\top B\mathbf{x}$,

$$\mathbf{x}^\top B\mathbf{x} = \mathbf{x}^\top W^\top AW\mathbf{x} = (W\mathbf{x})^\top A(W\mathbf{x}). \quad (13)$$

Because $W\mathbf{x} \neq \mathbf{0}$ and A is an SPD matrix, $\mathbf{x}^\top B\mathbf{x} > \mathbf{0}$. This proves that B is an SPD matrix.

On the other hand, if B is an SPD matrix, for arbitrary nonzero vector $\mathbf{x} \in \mathbb{R}^{C' \times 1}$, $\mathbf{x}^\top B\mathbf{x} = (W\mathbf{x})^\top A(W\mathbf{x}) > \mathbf{0}$. Because A is an SPD matrix, $W\mathbf{x} \neq \mathbf{0}$. Only if $\mathbf{x} = \mathbf{0}$ can lead to $W\mathbf{x} = \mathbf{0}$. Accordingly, $\text{rank}(W) = C'$ and W is a column full rank matrix. \square

Since there are learnable parameters in the SPD matrix transformation layer, we should not only compute the gradient of loss function \mathcal{L} with respect to the input K , but also calculate the gradient with respect to parameters W . The gradient with respect to the input K is

$$\frac{\partial \mathcal{L}}{\partial K} = W \frac{\partial \mathcal{L}}{\partial Y} W^\top, \quad (14)$$

where $\frac{\partial \mathcal{L}}{\partial Y}$ is the gradient with respect to the output Y .

Since W is a column full rank matrix, it is on a non-compact Stiefel manifold [38]. However, directly optimizing W on the non-compact Stiefel manifold is infeasible. To

overcome this issue, we relax W to be semi-orthogonal, i.e., $W^\top W = I_{C'}$. In this case, W is on the orthogonal Stiefel manifold $St(C', C)$. The optimization space of parameters W is changed from the non-compact Stiefel manifold to the orthogonal Stiefel manifold $St(C', C)$. Considering the manifold structure of W , the optimization process is quite different from the gradient descent method in the Euclidean space. We first compute the partial derivative with respect to W . Then we convert the partial derivative to the manifold gradient that lies on the tangent space. Along the tangent gradient, we find a new point on the tangent space. Finally, the retraction operation is applied to map the new point on the tangent space back to the orthogonal Stiefel manifold. Thus, an iteration of the optimization process on the manifold is completed. This process is illustrate in Fig. 4. Next we will elaborate each step.

First the partial derivative $\frac{\partial \mathcal{L}}{\partial W}$ with respect to W is computed by

$$\frac{\partial \mathcal{L}}{\partial W} = K^\top W \frac{\partial \mathcal{L}}{\partial Y} + KW \left(\frac{\partial \mathcal{L}}{\partial Y} \right)^\top. \quad (15)$$

The partial derivative $\frac{\partial \mathcal{L}}{\partial W}$ doesn't contain any manifold constraints. Considering W is a point on the orthogonal Stiefel manifold, the partial derivative needs to be converted to the manifold gradient, which is on the tangent space. As shown in Fig. 5, on the orthogonal Stiefel manifold, the partial derivative $\frac{\partial \mathcal{L}}{\partial W}$ is a Euclidean gradient at the point W , not tangent to the manifold. The tangential component of $\frac{\partial \mathcal{L}}{\partial W}$ is what we need for optimization, which lies on the tangent space. The normal component is perpendicular to the tangent space. We decompose $\frac{\partial \mathcal{L}}{\partial W}$ into two vectors that are perpendicular to each other, i.e., one is tangent to the manifold and the other is the normal component based on the Theorem 3.

Theorem 3. Let M denote an orthogonal Stiefel manifold and \mathbf{X} is a point on M . $F(\mathbf{X})$ denotes a function defined on the orthogonal Stiefel manifold. If the partial derivatives of F with respect to \mathbf{X} is $F_{\mathbf{X}}$, the manifold gradient ∇F at \mathbf{X} which is tangent to M is $\nabla F = F_{\mathbf{X}} - \mathbf{X}F_{\mathbf{X}}^\top \mathbf{X}$.

Proof. Because \mathbf{X} is a point on the orthogonal Stiefel manifold, $\mathbf{X}^\top \mathbf{X} = I$, where I is an identity matrix. Differentiating $\mathbf{X}^\top \mathbf{X} = I$ yields $\mathbf{X}^\top \Delta + \Delta^\top \mathbf{X} = 0$, where Δ is a tangent vector. Thus, $\mathbf{X}^\top \Delta$ is a skew-symmetric matrix. Note that, the canonical metric for the orthogonal Stiefel manifold at the point Y is $g_c(\Delta_1, \Delta_2) = \text{tr}(\Delta_1^\top (I - \frac{1}{2}YY^\top)\Delta_2)$. For all tangent vectors δ at \mathbf{X} , we can get that

$$\text{tr}F_{\mathbf{X}}^\top \Delta = g_c(\nabla F, \Delta) = \text{tr}(\nabla F)^\top (I - \frac{1}{2}\mathbf{X}\mathbf{X}^\top)\Delta. \quad (16)$$

Because $\mathbf{X}^\top (\nabla F)$ is a skew-symmetric matrix, Eq. (16) can be solved, i.e., $\nabla F = F_{\mathbf{X}} - \mathbf{X}F_{\mathbf{X}}^\top \mathbf{X}$. \square

Then the tangential component ∇L_W at W can be expressed by the partial derivative $\frac{\partial \mathcal{L}}{\partial W}$,

$$\nabla L_W = \frac{\partial \mathcal{L}}{\partial W} - W \left(\frac{\partial \mathcal{L}}{\partial W} \right)^\top W. \quad (17)$$

∇L_W is the manifold gradient of the orthogonal Stiefel manifold. Searching along the gradient ∇L_W gets a new point

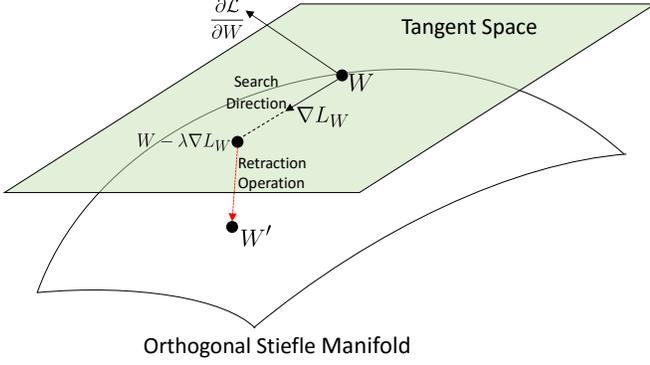


Fig. 4. The illustration of the optimization process of W . W is an original point on the orthogonal Stiefel manifold. W' is a new point after an iterative update. $\frac{\partial \mathcal{L}}{\partial W}$ is the partial derivative of the loss function with respect to W . ∇L_W is the manifold gradient lying on the tangent space.

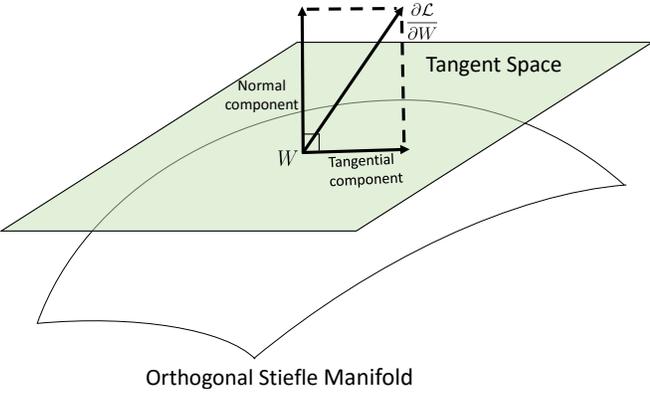


Fig. 5. The illustration of tangential and normal components of a vector. These two components are perpendicular. v is the vector on the surface. The tangential component of v is a vector on the tangent space and normal component is perpendicular to the tangent space.

on the tangent space. Finally, we use the retracting operation to map the point on the tangent space back to the Stiefel manifold space,

$$W := q(W - \lambda \nabla L_W), \quad (18)$$

where $q(\cdot)$ is the retraction operation mapping the data back to the manifold. Specifically, $q(A)$ denotes the Q matrix of QR decomposition to A . $A \in \mathbb{R}^{n \times p}$, $A = QR$, where $Q \in \mathbb{R}^{n \times p}$ is a semi-orthogonal matrix and $R \in \mathbb{R}^{p \times p}$ is an upper triangular matrix. λ is the learning rate.

Note that, we can make a Relu activation function layer follow the SPD matrix transformation layer. The output of the Relu layer is still an SPD matrix based on the Theorem 4.

Theorem 4. *The relu activation function on a matrix Y is $f(Y)$. Let $Z = f(Y)$,*

$$Z_{ij} = \begin{cases} 0, & \text{if } Y_{ij} < 0 \\ Y_{ij}, & \text{if } Y_{ij} \geq 0 \end{cases}.$$

If Y is an SPD matrix, Z is an SPD matrix.

Proof. The detailed proof of this theorem is shown in the appendix section of [24]. \square

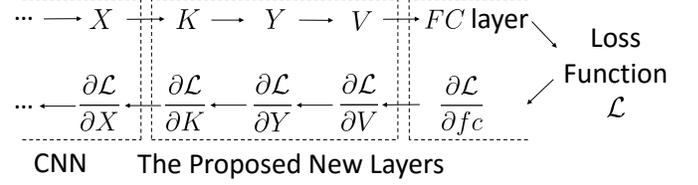


Fig. 6. Illustration of the forward and backward propagations of networks with the proposed aggregation method.

D. Vectorization layer

Since inputs of the common classifier is all vectors, we should vectorize the SPD matrix to a vector. Because of the symmetry of the robust SPD matrix achieved by the transformation layer, Y is determined by $\frac{C' \times (C'+1)}{2}$ elements, *i.e.*, the upper triangular matrix or the lower triangular matrix of Y . Here, we take the upper triangular matrix of Y and reshape it into a vector as the input of the loss function. Let's denote the vector by V ,

$$\begin{aligned} V &= [Y_{11}, \sqrt{2}Y_{12}, \dots, \sqrt{2}Y_{1C'}, Y_{22}, \\ &\sqrt{2}Y_{23}, \dots, \sqrt{2}Y_{C'(C'-1)}, Y_{C'C'}] \\ &= [V_1, V_2, \dots, V_{\frac{C' \times (C'+1)}{2}}]. \end{aligned} \quad (19)$$

Due to the symmetry of the matrix Y , the gradient $\frac{\partial \mathcal{L}}{\partial Y}$ is also a symmetric matrix. For the diagonal elements of Y , its gradient of the loss function is equal to the gradient of its corresponding element in the vector V , while the gradient of non-diagonal elements of Y is $\sqrt{2}$ times of the element in the vector V . The gradient with respect to Y is given by

$$\frac{\partial \mathcal{L}}{\partial Y} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial V_1} & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_2} & \dots & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{C'-1}} & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{C'}} \\ \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_2} & \frac{\partial \mathcal{L}}{\partial V_{C'+1}} & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{C'+2}} & \dots & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{2 \times C'-1}} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{C'}} & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{2 \times C'-1}} & \frac{\sqrt{2} \partial \mathcal{L}}{\partial V_{3 \times C'-3}} & \dots & \frac{\partial \mathcal{L}}{\partial V_{\frac{C' \times (C'+1)}{2}}} \end{bmatrix}. \quad (20)$$

The normalization operation is important as well. We use the power normalization ($V_i := \text{sign}(V_i) \sqrt{|V_i|}$) and l_2 normalization ($V := V/\|V\|_2$) operation following the vector V . The gradient formulation Eq. (9), Eq. (14) and Eq. (18) calculate the gradient with respect to the input of the corresponding layer, respectively. Once these gradients are obtained, the standard SGD backpropagation can be easily employed to update the parameters directly with the learning rate. Fig. 6 shows the data flow in our network with the proposed three layers including forward and backward propagations. fc denotes the output of the last fully-connected layer. In Algorithm 1, we summarize the training process of our model. We can use more than one SPD transformation layers in the network, where each one can be followed by a Relu layer as the activation layer.

TABLE I
COMPARISON FOR CNNs BASED METHODS IN TERMS OF AVERAGE PRECISION (%). OUR METHOD IS BOLD IN THE LAST LINE.

Method	DTD	FMD	KTH-T2b	CUB-200-2011	FGVC-aircraft
FV-CNN [31]	67.3	73.5	73.3	49.9	-
FV-FC-CNN [31]	69.8	76.4	73.8	54.9	-
B-CNN [8]	69.6	77.8	79.7	74.0	74.3
Deep-TEN _{ResNet50} [33]	-	80.2	82.0	-	-
VGG-16 [2]	66.8	77.8	78.3	68.0	75.0
Ours	68.9	79.2	81.1	72.4	77.8

Algorithm 1 Training Process of Our Model

Input: Training data $I = \{I_i\}_{i=1}^n$, where n is the number of training samples. SPD matrix transformation layer's initial parameters W . The other parameters in the CNN θ . Learning rate λ .

Output: SPD matrix transformation layer's parameters W . Other layer's parameters θ .

- 1: **while** not converge **do**
- 2: Compute the convolutional features M by forwarding I_i through convolutional layer.
- 3: Compute the SPD matrix K by Eq.(8).
- 4: Compute the transformed SPD matrix Y by Eq.(12).
- 5: Compute the vector representation V by Eq.(19).
- 6: Compute the loss \mathcal{L} .
- 7: Compute the gradient $\frac{\partial \mathcal{L}}{\partial V}$.
- 8: Compute the gradient $\frac{\partial \mathcal{L}}{\partial Y}$ by Eq.(20).
- 9: Compute the gradient $\frac{\partial \mathcal{L}}{\partial K}$ by Eq.(14).
- 10: Update the parameters W by Eq.(15), Eq.(17) and Eq.(18).
- 11: Compute the gradient $\frac{\partial \mathcal{L}}{\partial M}$ with respect to the last convolutional features.
- 12: Compute the gradient $\frac{\partial \mathcal{L}}{\partial \theta}$ with respect to the other parameters.
- 13: Update the other parameters $\theta = \theta - \lambda \frac{\partial \mathcal{L}}{\partial \theta}$.
- 14: **end while**
- 15: **return** θ and w

IV. EXPERIMENT

To demonstrate the benefits of our method, we conduct extensive experiments on visual classification tasks. We conduct experiments on visual classification tasks to show the performance of the SPD aggregation framework including the generation and transformation processes. We present the visual classification tasks on five datasets. We choose the challenging texture and fine-grained classification tasks. The texture classification tasks need a powerful global representation, because of the features of texture should be invariant to translation, scaling and rotation. Differences among fine-grained images are very small. It is challenging to represent these differences in the aggregation process.

A. Datasets and Evaluation Protocols

We choose three texture datasets in the experiments. They are Describable Textures Dataset (DTD) [39], Flickr Material Database (FMD) [40] and KTH-TIPS-2b (KTH-2b) [41]. DTD



Fig. 7. Example images on the DTD, FMD and KTH-2b datasets. We can see that images of the same category have huge differences, particularly the plastic class.



Fig. 8. Example reference images on the FGVC-aircraft and CUB-200-2011 datasets. The gaps between the pictures are very small.

and FMD are both collected in the wild conditions while KTH-2b is under the laboratory condition. DTD has 47 classes, and each class contains 120 images. There are totally 5640 images in DTD. FMD contains 1000 images of 10 classes, each class has 100 images. KTH-2b contains 4752 images of 11 classes. Fig. 7 illustrates the texture datasets for our experiments. For these texture datasets, we follow the standard train-test protocol. We divide DTD and FMD into three subsets randomly, and use two subsets for training and the rest one subset for testing. Images of KTH-2b are split into four samples. we train the framework using one sample and test

TABLE II
COMPARISON FOR THE COMPONENTS OF THE PROPOSED AGGREGATION METHOD IN TERMS OF AVERAGE PRECISION (%).

Method	DTD	FMD	KTH-T2b	CUB-200-2011	FGVC-aircraft
B-CNN [8]	67.9	77.8	79.7	74.0	74.3
512 conv 1×1 + B-CNN [8]	66.3	74.9	77.9	67.2	65.3
VGG-16 [2]	66.8	77.8	78.3	68.0	75.0
512 conv 1×1 + VGG-16 [2]	64.5	74.3	76.7	64.5	75.1
<i>no conv</i> 1×1 + Kernel Aggregation Layer	66.8	76.7	81.0	72.2	75.8
128 conv 1×1 + Kernel Aggregation Layer	67.8	78.3	81.1	64.6	72.3
512 conv 1×1 + Kernel Aggregation Layer	67.1	78.6	81.3	72.1	76.7
Ours	68.9	79.2	81.1	72.4	77.8

on the rest three samples. Inspired by [1], the texture images are augmented. We do 15 times augmentation to the training data, including randomly cropping 10 times and picking from center and four corners. The test images are only picked from the center and four corners. The size of cropped images are resized to 224×224 .

We report results on birds and aircrafts fine-grained recognition datasets. The birds dataset [42] is CUB-200-2011 which contains 200 classes and 11788 images totally. The FGVC-aircraft dataset [43] contains 10000 aircraft images of 100 classes. Fig. 8 illustrates some fine-grained images. We train and test the birds and aircrafts fine-grained datasets through the inherent training document. The data augmentation is not applied to the fine-grained images. We resize them to the size of 224×224 . All the texture and fine-grained images are normalized by subtracting means for RGB channels.

B. Implementation Details

The basic convolutional layers and pooling layers before our SPD aggregation are from the VGG-16 model which is pre-trained on the ImageNet dataset. We remove layers after the *conv5-3* layer of VGG-16 model. Then we insert our SPD aggregation method into the network following the *conv5-3* layer. Finally, a FC layer and a softmax layer follow the vectorization layer where the output dimension of the FC layer is equal to the number of classes. All our networks run under the caffe framework. We use SGD with a mini-batch size of 32. The training process is divided into two stages. At the first training stage, we fix the parameters before the SPD aggregation method and train the rest new layers. The learning rate is started from 0.1 and reduced by 10 when error plateaus. At the second training stage, we train the whole network. The learning rate is started from 0.001 and divided by 10 when error plateaus.

C. Experiments for the SPD Aggregation Framework

In this section, we compare the SPD aggregation framework with some state-of-the-art convolutional feature aggregation methods. First a 1×1 convolutional layer whose number of channels is 512 follows the *conv5-3* convolutional layers. Then our SPD aggregation method including a kernel aggregation layer, an SPD matrix transformation layer and a vectorization layer is inserted after the 1×1 convolutional layer. The output

size of the SPD matrix transformation layer is 512×512 . Considering these datasets are not big enough, we only use one SPD matrix transformation layer to avoid the overfitting. Table I shows the comparison on texture datasets and fine-grained datasets respectively.

The following methods are evaluated in our experiments, FV-CNN [31], FV-FC-CNN [31], B-CNN [8], Deep-TEN [33] and the pure VGG-16 model [2] is used as the baseline. FV-CNN aggregates convolutional features from VGG-16 *conv5-3* layer. The dimension of it is 65600 and is compressed to 4096 by PCA for classification. FV-FC-CNN incorporates the FC features and FV vector. B-CNN uses the Bilinear pooling method on the *conv5-3* layer of VGG-16 model. The Deep-TEN uses 50-layers ResNet and larger number of training samples and image size, while the other methods use VGG-16 model. It is not scientific to compare Deep-TEN with the other feature aggregation methods.

We can see that, our method gets a better performance, especially on KTH-2b and FGVC-aircraft datasets. The average precision of our method on KTH-2b and FGVC-aircraft datasets are 81.1% and 77.8%. In contrast, B-CNN achieves 79.7% and 74.3%. On DTD and CUB-200-2011 datasets, our method is slightly worse than the B-CNN. The reason may be that the linear relationships among features are dominant on some datasets and the nonlinear relationships are important on the others.

D. Experiments for the Components of the Proposed Aggregation Method

1×1 convolutional layer. As mentioned above, we employ a 1×1 convolutional layer to accomplish the preprocessing of convolutional features. In this section, we provide experiments for the necessity of the preprocessing of convolutional features in our method. We design experiments in Table II. We combine different numbers of channels of *conv* 1×1 layer with the kernel aggregation layer. We also add the *conv* 1×1 layer to the B-CNN and pure VGG network. *No conv* 1×1 , 128 *conv* 1×1 and 512 *conv* 1×1 in the Table indicate that whether there is a *conv* 1×1 layer before the kernel aggregation layer and the number of channels of the *conv* 1×1 layer. The kernel aggregation layer in the Table means that there is only the kernel aggregation layer without the SPD matrix transformation layer and the vectorization layer in the

network. Table II shows that, the *conv* 1×1 layer is beneficial to our nonlinear kernel aggregation method. However, it is useless or even harmful to the B-CNN and pure VGG-16 model. Our benefits are brought about by the powerful SPD matrix instead of the preprocessing of convolutional features. But the preprocessing of convolutional features can actually lead to better performance to the SPD aggregation. The reason may be that the convolutional features are totally different from the kernel matrices but have some similarities to the Bilinear matrices or FC features. We can also observe that the number of channels of *conv* 1×1 layer has small influence for the texture datasets. But when it is reduced to 128, the performance on fine-grained datasets is declined. So we argue that the convolutional features are redundant for the texture datasets but not redundant for the fine-grained datasets.

SPD Matrix Generation Process. To evaluate the effectiveness of the nonlinear SPD matrix generation process, we establish a network that only contains the kernel aggregation layer without the SPD matrix transformation layer and vectorization layer. The outcome is shown in Table II. Without the *conv* 1×1 layer, *i.e.*, no *conv* 1×1 + kernel aggregation layer in the Table, our network is comparable with the B-CNN and pure VGG-16 model. When the *conv* 1×1 layer is added to the network, *i.e.*, 512 *conv* 1×1 + Kernel Aggregation Layer in the Table, it has a obvious better performance than the other methods on FMD, KTH-2b and FGVC-aircraft datasets.

SPD Matrix Transformation Process. We design experiments to evaluate the effectiveness of the proposed transformation process in this subsection. Compared with ours, 512 *conv* 1×1 + kernel aggregation layer in Table II only lacks the SPD matrix transformation layer and the vectorization layer, the rest is the same. Through Table II, we find that the SPD matrix transformation process transforms the SPD matrix to a more suitable and discriminative representation. Especially on DTD dataset and FGVC-aircraft datasets, the performance is improved by 1.8% and 1.1% respectively.

V. CONCLUSION

In this paper, we have proposed a new powerful SPD aggregation method which models the convolutional feature aggregation as an SPD matrix non-linear learning problem on the Riemannian manifold. To achieve this goal, we have designed three new layers to aggregate the convolutional features into an SPD matrix and transform the SPD matrix to be more discriminative and suitable. The three layers include a kernel aggregation layer, an SPD matrix transformation layer and a vectorization layer under an end-to-end framework. We investigated the component decomposition and retraction of the Orthogonal Stiefle manifold to carry out the backpropagation of our model. Meanwhile, the faster matrix operation was adopted to speed up forward and backward backpropagations. Compared with alternative aggregation strategies such as FV, VLAD and bilinear pooling, our SPD aggregation achieves appealing performance on visual classification tasks. Extensive experiments on 11 challenging datasets have demonstrated that our approach outperforms the state-of-the-art methods.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems (NIPS)*, 2012, pp. 1097–1105.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1–9.
- [4] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [6] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2014, pp. 818–833.
- [7] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, "Decaf: A deep convolutional activation feature for generic visual recognition," in *Proceedings of the International Conference on Machine Learning (ICML)*, vol. 32, 2014, pp. 647–655.
- [8] T. Y. Lin, A. Roychowdhury, and S. Maji, "Bilinear cnn models for fine-grained visual recognition," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1449–1457.
- [9] A. B. Yandex and V. Lempitsky, "Aggregating local deep features for image retrieval," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2016, pp. 1269–1277.
- [10] G. Tolias, R. Sivic, and H. Jégou, "Particular object retrieval with integral max-pooling of cnn activations," *arXiv preprint arXiv:1511.05879*, 2015.
- [11] Y. Kalantidis, C. Mellina, and S. Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 685–701.
- [12] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell, "Compact bilinear pooling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 317–326.
- [13] Z. Huang, R. Wang, X. Li, W. Liu, S. Shan, L. V. Gool, and X. Chen, "Geometry-aware similarity learning on spd manifolds for visual recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. PP, no. 99, pp. 1–1, 2017.
- [14] M. Harandi, M. Salzmann, and R. Hartley, "Dimensionality reduction on spd manifolds: The emergence of geometry-aware methods," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 1–1, 2017.
- [15] Z. Huang, R. Wang, S. Shan, X. Li, and X. Chen, "Log-euclidean metric learning on symmetric positive definite manifold with application to image set classification," in *Proceedings of the International Conference on Machine Learning (ICML)*, vol. 33, 2015, pp. 720–729.
- [16] S. Herath, M. Harandi, and F. Porikli, "Learning an invariant hilbert space for domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3845–3854.
- [17] J. Zhang, L. Wang, L. Zhou, and W. Li, "Exploiting structure sparsity for covariance-based visual representation," *arXiv preprint arXiv:1610.08619*, 2016.
- [18] L. Zhou, L. Wang, J. Zhang, Y. Shi, and Y. Gao, "Revisiting metric learning for spd matrix based visual representation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3241–3249.
- [19] P. Li, J. Xie, Q. Wang, and W. Zuo, "Is second-order information helpful for large-scale visual recognition?" in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2070–2078.
- [20] C. Ionescu, O. Vantzos, and C. Sminchisescu, "Matrix backpropagation for deep networks with structured layers," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 2965–2973.
- [21] K. Yu and M. Salzmann, "Second-order convolutional neural networks," *Clinical Immunology and Immunopathology*, vol. 66, no. 3, pp. 230–238, 2017.
- [22] L. Wang, J. Zhang, L. Zhou, and C. Tang, "Beyond covariance: Feature representation with nonlinear kernel matrices," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 4570–4578.

- [23] L. Bo, X. Ren, and D. Fox, "Kernel descriptors for visual recognition," in *Advances in neural information processing systems (NIPS)*, 2010, pp. 244–252.
- [24] Z. Dong, S. Jia, C. Zhang, M. Pei, and Y. Wu, "Deep manifold learning of symmetric positive definite matrices with application to face recognition," in *Proceedings of the Association for the Advancement of Artificial Intelligence (AAAI)*, 2017, pp. 4009–4015.
- [25] Z. Huang and L. J. Van Gool, "A riemannian network for spd matrix learning," in *Proceedings of the Association for the Advancement of Artificial Intelligence (AAAI)*, 2017, pp. 2036–2042.
- [26] T. Zhang, W. Zheng, Z. Cui, Y. Zong, and Y. Li, "Deep manifold-to-manifold transforming network for action recognition," *arXiv preprint arXiv:1705.10732*, 2017.
- [27] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2003, pp. 1470–1477.
- [28] L. Liu, C. Shen, L. Wang, A. van den Hengel, and C. Wang, "Encoding high dimensional local features by sparse coding based fisher vectors," in *Advances in neural information processing systems (NIPS)*, 2014, pp. 1143–1151.
- [29] Y. H. Ng, F. Yang, and L. S. Davis, "Exploiting local features from deep networks for image retrieval," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2015, pp. 53–61.
- [30] Y. Gong, L. Wang, R. Guo, and S. Lazebnik, "Multi-scale orderless pooling of deep convolutional activation features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 392–407.
- [31] M. Cimpoi, S. Maji, and A. Vedaldi, "Deep filter banks for texture recognition and segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3828–3836.
- [32] L. Liu, C. Shen, and A. V. D. Hengel, "Cross-convolutional-layer pooling for image recognition," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 2305–2313, 2017.
- [33] H. Zhang, J. Xue, and K. Dana, "Deep ten: Texture encoding network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 708–717.
- [34] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5297–5307.
- [35] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [36] L. Z. X. L. Melih Engin, Lei Wang, "Deepkspd: Learning kernel-matrix-based spd representation for fine-grained image recognition," *arXiv preprint arXiv:1711.04047*, 2017.
- [37] K. Yu, B. Leng, Z. Zhang, D. Li, and K. Huang, "Weakly-supervised learning of mid-level features for pedestrian attribute recognition and localization," *arXiv preprint arXiv:1611.05603*, 2016.
- [38] P. A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*, 2009.
- [39] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi, "Describing textures in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 3606–3613.
- [40] L. Sharan, C. Liu, R. Rosenholtz, and E. H. Adelson, "Recognizing materials using perceptually inspired features," *International Journal of Computer Vision*, vol. 103, no. 3, pp. 348–371, 2013.
- [41] B. Caputo, E. Hayman, and P. Mallikarjuna, "Class-specific material categorisation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2005, pp. 1597–1604.
- [42] L. Xie, Q. Tian, R. Hong, S. Yan, and B. Zhang, "Hierarchical part matching for fine-grained visual categorization," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2013, pp. 1641–1648.
- [43] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi, "Fine-grained visual classification of aircraft," *arXiv preprint arXiv:1306.5151*, 2013.