



Sparse coding and normalization for deep Fisher score representation

Sixiang Xu, Damien Muselet, Alain Trémeau

Univ Lyon, UJM-Saint-Etienne, CNRS, Institut Optique Graduate School, Laboratoire Hubert Curien UMR 5516, F-42023, SAINT-ETIENNE, France

ABSTRACT

Fisher Scores have been shown to be accurate global image features for classification. However, their performance is very dependent on the quality of the input features as well as the normalization steps applied to them. In this paper, we propose to embed the Fisher vectors in an end-to-end trainable deep network by concentrating on these two crucial elements: adapting the encoding to the deep features and normalizing the extracted second order statistics. Therefore, we make use of a deep sparse coding module that allows to sample the center of each Gaussian function from a learned subspace and thus to better fit the high dimensional data distribution. Second, we introduce a new normalization module that computes an approximate square root matrix normalization well adapted to the Fisher vectors. These processing steps are embedded in a deep network so that all the modules work together for the sole purpose of improving classification performance. Experimental results show that this solution clearly outperforms many alternatives in the context of material, indoor scenes or fine-grained image classification.

© 2022 Elsevier Ltd. All rights reserved.

1. Introduction

Deep neural networks have emerged as an essential solution for performing classification tasks. In these networks, convolutional layers extract accurate local features that are pooled to a global feature vector which is sent to fully connected layers for classification. The first networks neglected the pooling step and directly sent the set of local features in the dense layers (Simonyan and Zisserman (2015)), while the series of ResNet apply a global average pooling to decrease the dimension of the global feature vector and hence reduce the number of parameters of the network (He et al. (2016)). Orderless pooling was widely used before convolutional neural networks (CNN) with the bags of visual words (BOW) (Lazebnik et al. (2006)), VLAD (Jégou et al. (2012)) or Fisher Vectors (Sánchez et al. (2013)) and has shown to provide good results when applied to CNN features (Cimpoi et al. (2015); Gong et al. (2014)). Among them, the Fisher Vectors (FV) were the most promising because they generalize the VLAD and BOW. The main idea of FV is to model the distribution of the training data with a Gaussian mixture and to characterize each data point with the derivatives over the model parameters. It appears that two main steps are crucial in such approach (Sánchez et al. (2013)): the

data distribution has to be accurately fitted by the Gaussian mixture and the provided second order statistics have to be carefully normalized. In this paper, we propose to embed the Fisher representation in an end-to-end trainable network by concentrating on these two steps.

First, a Gaussian Mixture Model (GMM) seems not to be well adapted to the deep local features since they are lying in a very high dimensional space and require too many Gaussians to be accurately modeled (Liu et al. (2014)). Liu et al. proposed a smart solution to overcome this problem which consists in sampling the center of each Gaussian from a subspace and therefore benefiting from an infinite number of Gaussians to fit the data distribution (Liu et al. (2014)). The authors show that this problem can be solved by a classical sparse coding method. Unfortunately, their approach can not take advantage of the main interest of the CNN, i.e. training end-to-end the feature extraction, the pooling and the classification layers. To cope with this problem, we propose in this paper, to make use of a deep sparse coding module proposed in (Gregor and LeCun (2010)).

Second, a recent study has shown that the normalization of the second order statistics has a strong impact on the classification performance (Lin and Maji (2017)). The authors proposed in particular to use a square-root matrix normalization combined with element-wise square-root and l_2 normalization for bi-linear pooling. Unfortunately, unlike the bi-linear pooling used in (Lin and Maji (2017)), our Fisher representation

e-mail:

sixiang.xu,damien.muselet,alain.tremeau@univ-st-etienne.fr
(Sixiang Xu, Damien Muselet, Alain Trémeau)

does not provide a square matrix, thus rendering the solution from (Lin and Maji (2017)) unusable. Thus, in this paper, we propose to adapt the square-root matrix normalization to non square matrices and to embed this original module in a deep network.

By combining these two main contributions, we propose an original end-to-end trainable deep network that extracts accurate feature from images, pools them into a deep Fisher representation and normalized these statistics. By backpropagating the gradient of the classification loss, we are able to make all these modules collaborate with the sole objective of improving the performance of the classification task. Experimental tests on three different datasets and three different backbone architectures show that our solution outperforms many alternatives.

This paper is an extended version of our previous works (Xu et al. (2021)), called hereafter E2E-SCF for end-to-end sparse coding Fisher vector, where only the sparse coding has been addressed. In this current version, we clearly improve the method, the results and the analysis over (Xu et al. (2021)) as follows:

- we address the problem of the normalization of the second order statistics,
- we force a zero-mean feature distribution for each image,
- we run extensive experiments with more network architectures and much more compared approaches,
- we propose a deep analysis of the impact of each module on the results with a clear ablation study,
- we provide the results of the original work (Xu et al. (2021)) and compare them with our results on three different datasets.

2. Related works

2.1. Orderless pooling

Orderless pooling was widely used before the emergence of the CNN-based solutions. The most popular approaches were based on bags of visual words (BOW) (Lazebnik et al. (2006)), VLAD (Jégou et al. (2012)) or Fisher Vectors (Sánchez et al. (2013)). Inspired by these early methods, some works have evaluated the Fisher vectors or VLAD from deep features for texture or image classification (Cimpoi et al. (2015); Gong et al. (2014)). They show improvements over the SIFT-based counterparts but, in their workflow, the dictionary or Gaussian mixture model are learned independently from the deep features and from the classifier, leaving a large margin of improvement.

Thus, the next works have focused on embedding orderless pooling in deep networks to allow end-to-end training. Passalis and Tefas have inserted a Bag-of-Features pooling in deep neural networks thanks to radial basis function neurons (Passalis and Tefas (2017)). The output of the pooling module is a histogram of the visual words (0^{th} order statistics) learned on the training set. And their variations also achieve remarkable performance in other tasks, such as color constancy (Laakom et al. (2020)), visual information analysis (Krestenitis et al. (2020)) and human action recognition (Yang et al. (2020))

Instead of counting the occurrences of the visual words in one image, VLAD-based approaches aggregate the residuals between the local features and their nearest visual words (1^{st} order statistics). NetVLAD (Arandjelović et al. (2016)) is the first network that implements VLAD and allows an end-to-end training for image retrieval task and then Deep Ten transforms VLAD as a residual module for image classification (Zhang et al. (2017)). Later, Xue et al. (2018), Hu et al. (2019) and Mao et al. (2021) improve Deep Ten in different aspects. It has been show that first order statistics are more accurate to characterize images in classification tasks and the Fisher vectors go further by using first and second order statistics. Deep FisherNet (Tang et al. (2019)) is an embedded implementation of the GMM Fisher vector. Lin et al. (2017) introduces NetFV which extends NetVLAD by appending the second order statistics. The end-to-end Fisher Vector is also applied in many different domains, like action recognition (Wang et al. (2019); Wang and Koniusz (2021)) and remote sensing image retrieval. The main disadvantage of all these approaches is that they rely on a limited number of codewords or Gaussian centers, which prevents accurate modeling of the data distribution in the high-dimensional deep feature spaces (Liu et al. (2014)).

One interesting solution to cope with this problem has been proposed by Li et al. (2017). The authors compute Fisher vectors from a mixture of factor analyzers (MFA), instead of the classical GMM. Their solution is embedded in a deep network which is trainable end-to-end. The idea of MFA is to approximate the data manifold by low dimensional linear spaces and, in this sense, is similar to the idea of sparse coding (Liu et al. (2014)). Nevertheless, even if the MFA module is embedded in a deep network, the authors show that an accurate initialization of the weights of the network is required to obtain good performance. This initialization consists in running an Expectation-Maximization algorithm on the set of local features that have to be saved in memory. Furthermore, it appears that this second order representation requires high computation costs, high number of parameters to learn and occupies a very large memory space (500k dimensions which is more than the image itself) (jacob et al. (2019)).

Another group of second-order pooling works is based on bilinear coding (Lin et al. (2017); Yu and Salzmann (2018); Yu et al. (2020, 2021)). For example, B-CNN is also an end-to-end trainable network and aggregates feature vectors by sum-pooling their outer products (Lin et al. (2017)). Since this pooled representation always has cumbersome size, the approaches SMSO, RUN and SRM propose to compress the bilinear pooled features and improves the classification performance (Yu and Salzmann (2018); Yu et al. (2020, 2021)). The results of these methods will be compared with ours in the experiments.

Our method is inspired by (Liu et al. (2014)), detailed in the next section. More recently, these authors have also proposed an improved version of their work in (Liu et al. (2017)), called HSCFV. It uses two dictionaries to code input features and consequently, doubles dimension size of the Fisher vector. Nevertheless, their approach is not embedded in a deep CNN for end-to-end training. Furthermore, as sparse coding module is a

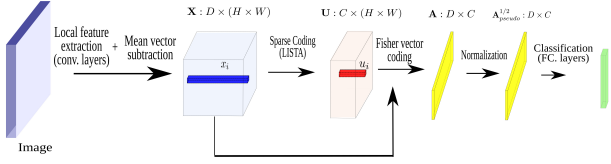


Fig. 1. Workflow of the proposed end-to-end trainable solution. First, deep features x_i are extracted with a classical CNN backbone and normalized (zero mean, Eq. 6). Then they are encoded into their sparse codes u_i with LISTA presented in Sec. 3.1.2. Next, Fisher score representation is produced (Eq. 5) and normalized with our proposed approach detailed in the Sec. 3.2.2.

core part in our approach, we investigate CNN-DL (Liu et al. (2018)) and SCN (Sun et al. (2019)) which propose their own implementation to produce sparse code. Unlike these methods using directly the sparse code as image features of first order statistics, our approach applies the code to further produce fisher score representation, which belongs to second order statistics and shows superior performance over CNN-DL (see Table 5)

Our method combines all the benefits of these previous solutions: it is embedded in an end-to-end trainable network, it samples an infinite number of Gaussian centers from a learned subspace and it does not require any heavy computation or storage to initialize the weights.

2.2. Normalization

As a post-processing step, after orderless pooling, normalization plays an important role in improving the performances. Perronnin et al. (2010) observed that the representation pooled by Fisher Vector is degraded by burstiness issues where discriminant but relatively rare visual features are overwhelmed by those that are more frequent. To alleviate this problem, some papers propose element-wise signed square rooting and L2-normalization (Perronnin et al. (2010); Arandjelovic and Zisserman (2013)). This normalization combination is also widely adopted in several successive orderless pooling works (Arandjelović et al. (2016); Lin et al. (2017); Liu et al. (2014)).

Besides the burstiness issue, Lin and Maji (2017) argued that the output of bilinear pooling should be normalized by matrix-logarithm functions in order to preserve the distances between elements in the manifold. Such normalization has been applied with success in (Carreira et al. (2012); Ionescu et al. (2015); Huang and Van Gool (2017)) with linear classifiers for semantic segmentation and image classification. The logarithm scales the eigenvalues in the Singular Value Decomposition (SVD) of a Symmetric Positive Definite (SPD) matrix A as $\log(A) = U \log(\Sigma) U^T$. Unfortunately, the SVD decomposition is computed inefficiently on GPUs (Lin and Maji (2017)), slowing down the network inference speed. Nevertheless, Lin and Maji propose a fast alternative approach with comparable performances and based on a variant of Newton iterations (Lin and Maji (2017)). This solution approximates the matrix square-root and can be embedded in a network that can be trained end-to-end.

Unfortunately, this approach is exclusively designed for SPD matrices such as the outputs of the bilinear pooling but can not

be directly applied to our Fisher representation that are rectangular and non symmetric matrices. Thus, we propose, in this paper, a new normalization step for such second order statistics matrices and that can also be embedded in a deep network.

3. Our approach

Fig. 1 illustrates the complete workflow of our solution whose successive steps are detailed in this Section. Our network starts with a pre-trained backbone constituted of convolutional layers, on top of which is applied an iterative sparse coding module called LISTA and detailed in Sec. 3.1.2. Then, the Fisher vectors are extracted from these features and normalized with our proposed solution. Then, dense layers provide the predicted categories.

3.1. Sparse Fisher coding

3.1.1. From subspace sampling to sparse coding

In order to increase the number of Gaussians that model the distribution of the data, we take advantage of the idea from (Liu et al. (2014)) that samples the Gaussian centers in a subspace spanned by a set of bases. Each mean vector is coded in this dictionary B with a code u drawn from a zero-mean Laplacian distribution (to enforce sparsity). Then each local feature vector x extracted from the images and associated with the code u is drawn from a Gaussian distribution $\mathcal{N}(Bu, \Sigma)$ centered on Bu . Fig. 2 illustrates the interest of this approach.

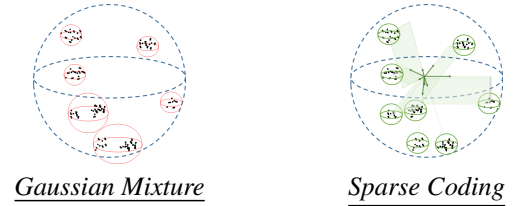


Fig. 2. Some data in a high dimensional space (illustrated by the sphere). Left: With GMM the data distribution is not well fitted because of the limited number of Gaussians. Right: With Sparse Coding, the Gaussian centers are coded sparsely in an adapted basis (green arrows) allowing to create unlimited number of Gaussians and so to fit better the data distribution. The sparsity is illustrated by the low number of basis required to code each center position (lines, planes or parallelograms).

Then, assuming a constant and diagonal covariance matrix as σ and using pointwise maximum to approximate the integral of the distribution, Liu et al. show that the logarithm of the likelihood of x can be estimated as (Liu et al. (2014)):

$$\log(P(x|B)) = \min_u \frac{1}{\sigma^2} \|x - Bu\|_2^2 + \lambda \|u\|_1, \quad (1)$$

where λ is the scale parameter of the Laplacian distribution of u .

Interestingly, this equation represents the classical problem of sparse coding. Liu et al. proposed to use an off-the-shelf sparse coding solver to learn the dictionary B and infer the code u . Obviously, making use of such independent solver is a good solution to minimize the reconstruction error of x with a sparse code, but it neglects the main goal which is to improve the performance of the classification task.

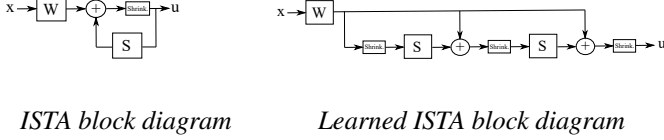


Fig. 3. Block diagrams of ISTA (left) and LISTA (right). ISTA evaluates the sparse code u of x with the iterative process detailed in Eq. 4. The shrinkage function is denoted \mathcal{T} in this equation. LISTA is an unfolded version of ISTA (2 iterations here) that can be embedded in an end-to-end trainable network. In our framework, the matrices S and W are learned and initialized thanks to our *warm-up* step, detailed in Sec.4.3.

Hence, we propose in the next section to embed a sparse coding module in a deep neural network that is trained end-to-end. The main advantage of such an approach is that it is learning a dictionary and sparse codes that are accurate to discriminate the different categories in the current dataset.

3.1.2. Embedding sparse coding with LISTA

Our aim is to find a solution for the following equation:

$$\min_u f(u) + \lambda \|u\|_1 \quad (2)$$

where $f(u) = \|x - Bu\|_2^2$, x is a data point, B the dictionary and u the sparse code of x .

One way to solve this equation is to resort to an Iterative Shrinkage/Thresholding Algorithm (ISTA) (Daubechies et al. (2004)) that iteratively approximates the solution with:

$$u_k = \mathcal{T}_{\lambda t_k}(u_{k-1} - t_k \nabla f(u_{k-1})), \quad (3)$$

where \mathcal{T}_α is a component-wise vector shrinkage function such that $[\mathcal{T}_\alpha(v)]_i = (|v_i| - \alpha)_+ \text{sign}(v_i)$, t_k is the step size at iteration k and ∇ is the gradient operator.

Evaluating the gradient of $f(u)$ defined above, we get:

$$\begin{aligned} u_k &= \mathcal{T}_{\lambda t_k}(u_{k-1} - 2t_k B^T (Bu_{k-1} - x)), \\ &= \mathcal{T}_{\lambda t_k}((I - 2t_k B^T B)u_{k-1} + 2t_k B^T x), \\ &= \mathcal{T}_{\lambda t_k}(S u_{k-1} + W x), \end{aligned} \quad (4)$$

where $S = I - 2t_k B^T B$ and $W = 2t_k B^T$.

As mentioned by Gregor and LeCun (2010), this equation can be illustrated as a recurrent block diagram as in Fig. 3, left. Fortunately, Gregor and LeCun (2010) proposed a fast approximation of ISTA called Learned ISTA (LISTA). This is an unfolded version of ISTA with a fix number of iterations and that can be plugged into a neural network to provide a sparse code (see Fig.3, right). Embedding this LISTA module in our CNN is a smart solution to learn a dictionary and sparse codes that help to discriminate between the categories of the current dataset.

3.1.3. Dictionary based Fisher coding

When a classical GMM is used to model the data distribution, the Fisher code is based on the partial derivatives of the posterior probabilities with respect to the weights, the mean and the standard-deviation parameters of the model (Sánchez et al. (2013)). In our case, inspired by Liu et al. (2014), we

use a particular Fisher coding, evaluated as the partial derivative of the log probability of the local features with respect to the dictionary itself:

$$\frac{\partial \log(P(x|B))}{\partial B} = \frac{\partial \frac{1}{\sigma^2} \|x - Bu^*\|_2^2 + \lambda \|u^*\|_1}{\partial B} = (x - Bu^*)u^{*T}, \quad (5)$$

where $u^* = \text{argmax}_u P(x|u, B)P(u)$ (see Liu et al. (2014) for details).

This module is very easy to insert in our deep network and provides the pooled features from the input image. These features are then sent to the last fully connected layers for classification. All these modules are constituting our CNN which can be trained end-to-end (see Fig. 1).

3.1.4. Mean Vector Subtraction

It is worth mentioning that, for each image, the input local feature vectors x_i are centered to have a zero mean before applying this dictionary encoding:

$$x'_i = x_i - \frac{1}{HW} \left(\sum_{i=1}^{HW} x_i \right), \quad (6)$$

where HW is the number of feature vectors x_i .

This pre-processing is similar to an instance normalization Ulyanov et al. (2016) without additive parameters to learn. It improves the generalization property of the model as observed in the experimental results (see more details in Sec. ??).

3.2. Fisher vector normalization

As mentioned earlier, the second order statistics tend to excessively emphasize very few coordinates, ignoring potential discriminant features (Lin and Maji (2017)). To cope with this problem, many normalization solutions have been proposed. In this paper, we take advantage of the approach proposed in (Lin and Maji (2017)) to normalize our Fisher vectors. Below, we first detail the solution of Lin and Maji (2017) and then, explain its extension to non square matrices.

3.2.1. Bilinear square matrix normalization

Assuming that the network backbone provides a feature map $X \in \mathbb{R}^{D \times H \times W}$ (see Fig.1), where H, W, D are the height, width and depth. This set of local feature vectors can be orderless pooled into a global feature vector by using bilinear pooling (Gao et al. (2016)). Therefore, the feature map X is reshaped to a 2D matrix ($D \times HW$) where each column x_i is ($D \times 1$) a local feature vector. Then, the output of the bilinear pooling is evaluated as:

$$A = \frac{1}{HW} \left(\sum_{i=1}^{HW} x_i x_i^T \right). \quad (7)$$

A is a ($D \times D$) symmetric positive definite (SPD) matrix.

While element-wise square-root normalization helps in improving the performance of the complete framework, Lin and Maji have shown that the results can be further boosted by applying a spectral normalization, i.e. scaling the eigenvalues of the associated covariance matrix (Lin and Maji (2017)). One way to do that is to transform the matrix A to its square-root

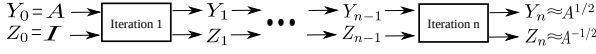


Fig. 4. Square-root matrix estimation with the Newton method. The inputs are a symmetric positive definite (SPD) matrix A and the identity I . After n iterations the outputs Y_n and Z_n converge into the square root matrix $A^{1/2}$ and the inverse square root matrix $A^{-1/2}$ of the input matrix A . One iteration is detailed in Fig.5.

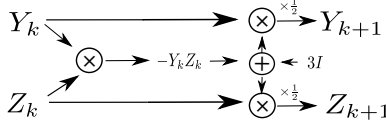


Fig. 5. One iteration of the Newton method as presented in Eq. 8.

$A^{1/2} = U\Sigma^{1/2}U^T$, where $A = U\Sigma U^T$ is the singular value decomposition (SVD) of A .

However, the computation of the SVD is poorly supported on GPUs and the authors suggest applying a variant of the Newton method to solve $F(Z) = Z^2 - A = 0$. Their approach is an iterative process where each iteration is as follow:

$$\begin{aligned} Y_{k+1} &= \frac{1}{2}Y_k(3I - Z_kY_k), \\ Z_{k+1} &= \frac{1}{2}(3I - Y_kZ_k)Z_k. \end{aligned} \quad (8)$$

By initializing $Y_0 = A$ and $Z_0 = I$, Y_k and Z_k converge to $A^{1/2}$ and $A^{-1/2}$ in very few iterations (even one) and requires only matrix multiplications (no inverse). The process is illustrated in Fig. 4 and 5.

This matrix normalization clearly improves the accuracy and efficiency of the bilinear pooling CNN (Lin and Maji (2017)), but it can't be directly applied to our Fisher representation, as shown in the next section.

3.2.2. Matrix Normalization for Fisher score representation

As previously explained, our Fisher representation is also a second-order matrix that could benefit from spectral normalization. From eq. 5, we know that it is expressed as:

$$A = \frac{1}{HW} \left(\sum_{i=1}^{HW} (x'_i - Bu_i)u_i^T \right), \quad (9)$$

where $B \in \mathbb{R}^{D \times C}$ is a dictionary (with C codewords) and $u_i \in \mathbb{R}^{C \times 1}$ is the sparse code of x'_i .

This matrix $A \in \mathbb{R}^{D \times C}$ is neither square nor symmetric and thus, can't be used as input for the Newton normalization that is restricted to SPD matrices. Indeed, since A is not SPD, its SVD is given as $A = U\Sigma V^T$, where $U \neq V$ and where $\Sigma \in \mathbb{R}^{D \times C}$ is not square.

In order to apply spectral normalization, we propose to estimate a so-called pseudo square root matrix $A_{pseudo}^{1/2}$ defined as:

$$A_{pseudo}^{1/2} = U\Sigma_{pseudo}^{1/2}V^T, \quad (10)$$

where $\Sigma_{pseudo}^{1/2}$ is calculated by square rooting the diagonal elements of Σ . Note that there is no matrix $\Sigma^{1/2}$ such that $\Sigma = \Sigma^{1/2}\Sigma^{1/2}$.

Inspired by (Lin and Maji (2017)), in order to avoid SVD computation, we resort to the Newton method to evaluate such a $A_{pseudo}^{1/2}$ matrix. But, since this solution only accepts SPD square matrix as input, we transform A into a square SPD matrix D evaluated as:

$$D = A^T A = V\Sigma^T U^T U \Sigma V^T = V\Sigma^T \Sigma V^T. \quad (11)$$

Note that this transform does not depend on U .

Since Σ is not symmetric, we introduce a helper matrix $H = [I_C | 0]^T \in \mathbb{R}^{D \times C}$, with I_C the $C \times C$ identity matrix, such that Σ can be expressed as:

$$\Sigma = H\tilde{\Sigma}, \quad (12)$$

where $\tilde{\Sigma}$ is a $C \times C$ square diagonal matrix.

Hence, Eq.11 can be derived into:

$$D = V\Sigma^T \Sigma V^T = V\tilde{\Sigma}^T H^T H \tilde{\Sigma} V^T = V\tilde{\Sigma}^2 V^T. \quad (13)$$

This equation is the SVD of the matrix D .

Feeding the previous Newton workflow with D and an identity matrix, we obtain $D^{1/2} = V\tilde{\Sigma}V^T$ and $D^{-1/2} = V\tilde{\Sigma}^{-1}V^T$ and feeding again this workflow with $D^{1/2}$ and an identity matrix, we obtain $D^{1/4} = V\tilde{\Sigma}^{1/2}V^T$ and $D^{-1/4} = V\tilde{\Sigma}^{-1/2}V^T$.

Finally, we have access to $A_{pseudo}^{1/2}$ thanks to:

$$\begin{aligned} AD^{-1/4} &= U\Sigma V^T V\tilde{\Sigma}^{-1/2}V^T, \\ &= U\Sigma\tilde{\Sigma}^{-1/2}V^T, \\ &= UH\tilde{\Sigma}\tilde{\Sigma}^{-1/2}V^T, \\ &= UH\tilde{\Sigma}^{1/2}V^T, \\ &= U\Sigma_{pseudo}^{1/2}V^T, \\ &= A_{pseudo}^{1/2}. \end{aligned} \quad (14)$$

Hence, without any SVD computation, this solution allows us to spectrally normalize a non SPD matrix A as $A_{pseudo}^{1/2}$ very efficiently. Furthermore, this workflow can be easily embedded in an end-to-end trainable deep network.

Thus, we propose to apply this new spectral normalization to our sparse Fisher encoding for classification tasks. All the framework can be trained end-to-end. In the next section, we propose to run extensive tests on different datasets to assess the quality of this method.

4. Experiments

In order to show that our solution generally helps the classification performance, we run experiments on three datasets, which vary between tasks and scales. The three datasets and their experimental settings are detailed in Section 4.1 and 4.2. Next, the training strategy of our network is shown in Section 4.3. In Section 4.4, the results and comparisons will be discussed.

4.1. Datasets

Orderless pooling methods were originally designed for texture and material recognition tasks (Lin et al. (2017); Yu and

Salzmann (2018); Zhang et al. (2017)). So we have first selected a reference material dataset. Then, these approaches have also been shown to provide good results on scene classification as well as fine-grained image classification (Liu et al. (2017); Lin et al. (2017); Li et al. (2017); Yu et al. (2021)). Thus, we have also selected two dedicated datasets for these tasks. The choice of these three datasets (detailed hereafter) is also a good way to validate the versatility of our solution for different image classification tasks.

The dataset MINC-2500 (Bell et al. (2015)), containing 23 commonly-seen material categories and 2,500 images per category, is a challenging large-scale dataset as material shows great intra-class variability in the real-world environment. The dataset MIT Indoor 67 (Quattoni and Torralba (2009)) is a medium but widely accepted benchmark for indoor scene classification task with 67 indoor categories and 100 images in each category. The dataset CUB-200-2011 (Wah et al. (2011)) provides 11,788 images of 200 bird species and is considered as a fine-grained classification dataset because the inter-class differences between bird species are subtle and sometime barely noticeable. In our experiments, we don't use the available object bounding boxes and part annotations. Note that we always make use of official training-test splits released with the datasets.

4.2. Experimental settings

Deep Pooling Module (DPM) - Our DPM is composed of a 1×1 convolution layer, a LISTA module with two iterations (see Fig. 3), the Fisher encoding layer and normalization process which includes matrix normalization (see section 3.1), element-wise square root and l_2 normalization. Then, the DPM is followed by a fully connected layer with softmax activation for classification.

Depending on dataset scales and for fair comparison with other works, we use different backbones and training strategies.

MIT-67 and CUB-200 settings - We adopt the settings of the state-of-the-art (Yu and Salzmann (2018); Lin et al. (2017)). The input image size is 448x448 and the backbone networks are either the pretrained VGG-D (a.k.a VGG-16) or Alexnet. Our DPM is plugged after the ReLU activation of the last convolutional layer. The 1×1 convolutional layer in the DPM does not change the input feature size and the sparse code in LISTA has 100 elements.

MINC-2500 settings - The network backbone is the pretrained ResNet-50 (He et al. (2016)). With the 1×1 convolutional layer in the DPM, the input feature size is reduced to 128 and the size of sparse code in LISTA is 32. While training, we follow the data augmentation settings from (Xue et al. (2018)). First, the input image is resized to 256×256 . Then we crop each image at i) a random location with ii) a random size (between 8% to 100% of the image area) and iii) a random aspect ratio (between 3/4 and 4/3). The crop is resized to 224×224 and used as the network input.

4.3. Training details

In the training phase, three consecutive steps are conducted. First, we run a PCA on a small subset of feature vectors (around 10,000) extracted from the backbone outputs and initialize the

Table 1. Ablation study of our workflow on the MIT-67 dataset. Essential elements in our approach are progressively added and the accuracy(%) given by their different combination is measured, showing their individual contribution to the classification.

LISTA	Warm-Up	Mean Sub.	Matrix Norm.	Accuracy
				76.72
✓				77.16
✓	✓			80.22
✓	✓	✓		80.60
	✓	✓	✓	80.67
✓	✓	✓	✓	81.24

1×1 convolutional layer of our DPM with these PCA parameters. Second, inspired by Branson et al. (2014), we apply a warm-up process that consists in training our DPM and FC layer (while the backbone is frozen) with an objective function which is the sum of the cross-entropy loss and the sparse coding loss (see Eq. (1)). Finally, the whole network is fine-tuned end-to-end under the supervision of the sole cross-entropy loss.

The optimization algorithm is a gradient descent with a mini-batch size of 64, a weight decay of $5e^{-4}$ and a momentum of 0.9. The learning rate is 0.004 during the warm-up. During the end-to-end finetuning, it starts from 0.004 and is divided by 10 when the training loss meets a plateau.

4.4. Results

In this Section, we provide many results in order to assess the quality of each contribution, to measure the impact of the hyperparameters and to compare our whole framework with the state-of-the-art.

Ablation study - In order to measure the impact of each of our different contributions, we propose to conduct an ablation study. The tests are run on the MIT-67 dataset with the VGG-16 network and the results are provided in Table 1.

For this study, we propose to start from the baseline network without contributions and to consecutively add the proposed modules in order to assess their individual impact on the results. When the LISTA module is not in the network, it is replaced by a 1×1 convolutional layer providing the codes u_i .

As introduced in Sec 4.3, our warm-up process is one of the three steps in the training phase. The goal is to train our DPM and the FC layer before fine-tuning the whole network. We can see in Table 1, that this training step boosts the performance from 77.16% to 80.22%, showing that an accurate initialization is important for our DPM and classifier.

Likewise, we notice that the proposed matrix normalization (called *Matrix Norm.* in Table 1) is one key element of our framework since it improves the accuracy from 80.60% to 81.24%.

Furthermore, centering the deep features (*Mean Sub.* in Table 1) thanks to eq. 6 also provides a slight improvement from 80.22% to 80.60%.

Finally, the impact of the LISTA module is measured with two different tests. Starting from the baseline and adding LISTA improves the results from 76.72% to 77.16% and adding

LISTA to the whole process helps to increase from 80.67% to 81.24%.

This ablation study is also a nice way to measure the improvement of our contributions over our previous work called E2E-SCF (Xu et al. (2021)). Indeed, in Table 1, the row with LISTA and Warm-up corresponds to our E2E-SCF. We notice that the additive contributions help to improve the accuracy from 80.22% to 81.24% on this dataset. Additional comparisons with this previous paper are proposed in Table 5 with 3 architectures and 3 datasets.

After analyzing the contribution of each element, we propose to discuss their individual computational costs. Thus, we have measured their inference times on the MIT-67 dataset with the VGG-16 backbone. According to Table 2, the feature extraction with the convolutional backbone (VGG-16, here) is clearly the bottleneck of the framework. The inference times of our proposed blocks are negligible compared to the one of the backbone. Among our proposed modules, the matrix normalization, i.e. the Newton algorithm, has the highest computational cost.

Table 2. Inference time (ms) per mini-batch (64 samples) required by each element of our framework. *Backb.* represents the convolutional layers used to extract the deep features, *Norm.* is our matrix normalization step and *Fisher* is the Fisher score encoding.

	Backb.	Mean Sub.	LISTA	Fisher	Norm.	Classif.
Time	1493	0.7	5.3	3.7	30.9	0.8

Hyperparameters - In order to go a step further in the analysis of our framework, we propose to study the impacts of two hyper-parameters on the results; namely the number of iterations in the LISTA module and the size of the dictionary in the sparse coding. Like the previous experiment, the tests are conducted on the MIT-67 dataset with the VGG-16 network.

LISTA is an unfolded version of ISTA and the number of iterations is an hyper-parameter. We investigate the performance of our framework across different numbers of iterations from 0 to 5, where 0 means that the LISTA module is replaced by a 1×1 convolutional layer. In Table 3, we notice that 2 or 3 iterations provide the best performance. After 3 iterations, the results start decreasing. Our intuition is that too many iterations of LISTA produce sparser codes at the expense of classification accuracy. For all the other tests in this paper, 2 iterations are used.

We also conducted an analysis on the number of codewords (dictionary size) required in the LISTA module. We measure the classification accuracy for a range of codeword numbers from 50 to 512 in Table 4. We notice that, due to overfitting, when the number of codewords is higher than 100, lower accuracy is observed. The default value for the next tests is 100.

Comparison with state-of-the-art - The top-1 classification

Table 3. Impact of the number of iterations in LISTA on the accuracy(%).

Iter. number	0	1	2	3	4	5
Accuracy	80.67	81.04	81.24	81.34	81.04	80.30

Table 4. Impact of the dictionary size on the accuracy (%).

Dic. size	50	100	200	300	400	512
Accuracy	80.37	81.24	80.90	80.00	80.22	80.22

accuracy of our approach and many alternatives are provided in Table 5. The results of the related works are directly extracted from the reference papers cited in the Table. Note that our CNN is trained on single-scale images while many state-of-the-art approaches are trained on multi-scales, so we have carefully selected the results that allows fair comparisons, but still some results in Table 5 are from multi-scale training (see comments in Table 5).

The methods called *Off-the-shelf* use independent modules that are not fine-tuned together while the *Finetuned* group contains approaches that use end-to-end trainable networks. We notice that the results provided by fine-tuned networks overall outperform those of the Off-the-shelf solutions. This shows that it is better to make the modules work together to optimize the same loss instead of independently optimizing them. Besides end-to-end learning attribute, our approach is built upon Deep Fisher Score Representation via Sparse Coding (SCFVC Liu et al. (2014)) and our E2E-SCF (Xu et al. (2021)) which produce more discriminant second-order pooled features than the classical Fisher vector or VLAD. We can see in Table 5 that the proposed smart combination of these two advantages make our method outperform the alternatives for all the tested datasets and backbones.

5. Conclusion

In this paper, we have proposed a complete workflow to extract second order statistics from images in the context of image classification. The approach is based on Fisher encoding which requires a data distribution fitting with Gaussians. We have first proposed to sparsely encode the Gaussian centers in a learned basis in order to improve the data fitting. Second, since the second order features require a spectral normalization before being used for classification, we have introduced an original matrix normalization based on a Newton algorithm. The main advantage of these two modules is that they can be embedded in a deep network that can be trained end-to-end. We have also proposed a training strategy that can easily initialize the network parameters before finetuning. Many experimental tests clearly show that our method outperforms the recent alternatives on three different datasets. The proposed non-SPD matrix normalization can be exploited to improve other second order statistics features such as those provided by compact bilinear coding. This is the aim of our future works.

References

- Arandjelović, R., Gronat, P., Torii, A., Pajdla, T., Sivic, J., 2016. Netvlad: Cnn architecture for weakly supervised place recognition, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Arandjelovic, R., Zisserman, A., 2013. All about vlad, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pp. 1578–1585.

Table 5. Comparison of the classification accuracy (%) with closed-related alternatives on three datasets and three backbone architectures.

	Approaches	MIT AlexNet	MIT VGG16	MIT ResNet50	CUB AlexNet	CUB VGG16	CUB ResNet50	MINC ResNet50
Off-shelf	Baseline (Sharif Razavian et al. (2014); Lin et al. (2017))	58.4			53.3	60.4		
	GMMFVC (Liu et al. (2014, 2017))	64.3	72.6 ^a		61.7	70.1 ^a		
	SCFVC (Liu et al. (2014, 2017))	68.2	77.6 ^a	83.28	66.4	77.3 ^a	78.86	
	HSCFVC (Liu et al. (2017))		79.5 ^a			80.8		
End-to-end	Baseline (Lin et al. (2017); Yu and Salzmann (2018))		64.51	76.45		70.4	74.51	79.1
	Deep Ten (Zhang et al. (2017))			71.3				80.4
	NetVLAD (Lin et al. (2017))					81.9		
	NetFV (Lin et al. (2017))		78.2			79.9		
	FisherNet (Tang et al. (2019))		76.4					
	MFAFVNet (Li et al. (2017))	69.89 ^b	78.01 ^b					
	CNN-DL (Liu et al. (2018))	66.60	78.33					
	B-CNN (Lin et al. (2017); Yu and Salzmann (2018))		77.6			84.0	79.05	
	SMSO (Yu and Salzmann (2018))		79.45	79.68		85.01	85.77	81.3
	SRM (Yu et al. (2021))		80.3			85.5		
	RUN (Yu et al. (2020))		80.5			85.7		
	E2E-SCF (Xu et al. (2021))	70.15	80.22	84.85	76.8	84.28	84.47	81.5
	Ours	70.60	81.24	85.52	77.49	85.8	87.38	81.8

^a These methods were trained with VGG19 (not VGG16) with 2 scales, whereas the other approaches from the column are trained with a single scale.

^b Since MFAFVNet works on patches and not on images, we have selected in Li et al. (2017) the results provided with the nearest patch scale from our settings (160 × 160).

- Bell, S., Upchurch, P., Snavely, N., Bala, K., 2015. Material recognition in the wild with the materials in context database. *Computer Vision and Pattern Recognition (CVPR)*.
- Branson, S., Van Horn, G., Belongie, S., Perona, P., 2014. Bird species categorization using pose normalized deep convolutional nets. *arXiv preprint arXiv:1406.2952*.
- Carreira, J., Caseiro, R., Batista, J., Sminchisescu, C., 2012. Semantic segmentation with second-order pooling, in: *European Conference on Computer Vision*, Springer. pp. 430–443.
- Cimpoi, M., Maji, S., Vedaldi, A., 2015. Deep filter banks for texture recognition and segmentation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3828–3836.
- Daubechies, I., Defrise, M., Mol, C., 2004. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Comm. Pure Appl. Math.* 57, 1413–1457.
- Gao, Y., Beijbom, O., Zhang, N., Darrell, T., 2016. Compact bilinear pooling, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 317–326.
- Gong, Y., Wang, L., Guo, R., Lazebnik, S., 2014. Multi-scale orderless pooling of deep convolutional activation features, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 392–407.
- Gregor, K., LeCun, Y., 2010. Learning fast approximations of sparse coding, in: *Proc. International Conference on Machine learning (ICML'10)*.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hu, Y., Long, Z., AlRegib, G., 2019. Multi-level texture encoding and representation (multer) based on deep neural networks, in: *2019 IEEE International Conference on Image Processing (ICIP)*, IEEE. pp. 4410–4414.
- Huang, Z., Van Gool, L., 2017. A riemannian network for spd matrix learning, in: *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Ionescu, C., Vantzos, O., Sminchisescu, C., 2015. Matrix backpropagation for deep networks with structured layers, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2965–2973.
- Jacob, P., Picard, D., Histace, A., Klein, E., 2019. Efficient codebook and factorization for second order representation learning, in: *Proc. International Conference on Learning Representations (ICLR)*.
- Jégou, H., Perronnin, F., Douze, M., Sánchez, J., Pérez, P., Schmid, C., 2012. Aggregating local image descriptors into compact codes. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 34(9).
- Krestenitis, M., Passalis, N., Iosifidis, A., Gabbouj, M., Tefas, A., 2020. Re-current bag-of-features for visual information analysis. *Pattern recognition* 106, 107380.
- Laakom, F., Passalis, N., Raitoharju, J., Nikkanen, J., Tefas, A., Iosifidis, A., Gabbouj, M., 2020. Bag of color features for color constancy. *IEEE Transactions on Image Processing* 29, 7722–7734.
- Lazebnik, S., Schmid, C., Ponce, J., 2006. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories, in: *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, pp. 2169–2178.
- Li, Y., Dixit, M., Vasconcelos, N., 2017. Deep scene image classification with the mfaivnet, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5746–5754.
- Lin, T.Y., Maji, S., 2017. Improved bilinear pooling with cnns, in: *Tae-Kyun Kim, Stefanos Zafeiriou, G.B., Mikolajczyk, K. (Eds.), Proceedings of the British Machine Vision Conference (BMVC)*, BMVA Press. pp. 117.1–117.12. URL: <https://dx.doi.org/10.5244/C.31.117>, doi:10.5244/C.31.117.
- Lin, T.Y., RoyChowdhury, A., Maji, S., 2017. Bilinear convolutional neural networks for fine-grained visual recognition. *IEEE transactions on pattern analysis and machine intelligence* 40, 1309–1322.
- Liu, L., Shen, C., Wang, L., Hengel, A.v.d., Wang, C., 2014. Encoding high dimensional local features by sparse coding based fisher vectors, in: *Advances in Neural Information Processing Systems (NIPS)*.
- Liu, L., Wang, P., Shen, C., Wang, L., Van Den Hengel, A., Wang, C., Shen, H.T., 2017. Compositional model based fisher vector coding for image classification. *IEEE transactions on pattern analysis and machine intelligence* 39, 2335–2348.
- Liu, Y., Chen, Q., Chen, W., Wassell, I., 2018. Dictionary learning inspired deep network for scene recognition, in: *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Mao, S., Rajan, D., Chia, L.T., 2021. Deep residual pooling network for texture recognition. *Pattern Recognition* 112, 107817.
- Passalis, N., Tefas, A., 2017. Learning bag-of-features pooling for deep convolutional neural networks, in: *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 5766–5774.
- Perronnin, F., Sánchez, J., Mensink, T., 2010. Improving the fisher kernel for large-scale image classification, in: *European conference on computer vision*, Springer. pp. 143–156.
- Quattoni, A., Torralba, A., 2009. Recognizing indoor scenes, in: *2009 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE. pp. 413–420.

- Sharif Razavian, A., Azizpour, H., Sullivan, J., Carlsson, S., 2014. Cnn features off-the-shelf: An astounding baseline for recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops.
- Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, in: Proc. International Conference on Learning Representations (ICLR'15).
- Sun, X., Nasrabadi, N.M., Tran, T.D., 2019. Supervised deep sparse coding networks for image classification. *IEEE Transactions on Image Processing* 29, 405–418.
- Sánchez, J., Mensink, T., Verbeek, J., 2013. Image classification with the fisher vector: Theory and practice. *International Journal of Computer Vision* 105.
- Tang, P., Wang, X., Shi, B., Bai, X., Liu, W., Tu, Z., 2019. Deep fishnet for image classification. *IEEE transactions on neural networks and learning systems* 30, 2244–2250.
- Ulyanov, D., Vedaldi, A., Lempitsky, V., 2016. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*.
- Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S., 2011. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001. California Institute of Technology.
- Wang, L., Koniusz, P., 2021. Self-supervising action recognition by statistical moment and subspace descriptors, in: Proceedings of the 29th ACM International Conference on Multimedia, pp. 4324–4333.
- Wang, L., Koniusz, P., Huynh, D.Q., 2019. Hallucinating idt descriptors and i3d optical flow features for action recognition with cnns, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 8698–8708.
- Xu, S., Muselet, D., Trémeau, A., 2021. Deep fisher score representation via sparse coding, in: Proceedings of the 19th International Conference on Computer Analysis of Images and Patterns (CAIP), Springer Verlag's series Lecture Notes in Computer Science.
- Xue, J., Zhang, H., Dana, K., 2018. Deep texture manifold for ground terrain recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 558–567.
- Yang, J., Liu, W., Yuan, J., Mei, T., 2020. Hierarchical soft quantization for skeleton-based human action recognition. *IEEE Transactions on Multimedia* 23, 883–898.
- Yu, K., Salzmann, M., 2018. Statistically-motivated second-order pooling, in: Proceedings of the European Conference on Computer Vision (ECCV), pp. 600–616.
- Yu, T., Cai, Y., Li, P., 2020. Toward faster and simpler matrix normalization via rank-1 update, in: European Conference on Computer Vision, Springer. pp. 203–219.
- Yu, T., Li, X., Li, P., 2021. Fast and compact bilinear pooling by shifted random maclaurin, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 3243–3251.
- Zhang, H., Xue, J., Dana, K., 2017. Deep ten: Texture encoding network, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 708–717.