

# Learning of Separable Filters by Stacked Fisher Convolutional Autoencoders

Arash Shahriari  
arash.shahriari@anu.edu.au;csiro.au

Australian National University (ANU)  
Commonwealth Scientific & Industrial  
Research Organisation (CSIRO)  
Canberra, Australia

## Abstract

Learning of convolutional filters in deep neural networks proves high efficiency to provide sparse representations for the purpose of image recognition. The computational cost of these networks can be alleviated by focusing on separable filters to reduce the number of learning parameters. Autoencoders are a family of powerful deep networks to build scalable generative models for automatic feature learning. Inspired by their stacked hierarchy, we introduce Fisher convolutional autoencoders to learn separable filters in a distributed architecture. These novel overcomplete autoencoders employ discriminant analysis to impose the highest possible distinction among texture classes whilst holds the minimum separation within each individual class. A distributed network of stacked Fisher autoencoders learns banks of separable filters in parallel and makes an ensemble of deep convolutional features with higher separability for a better classification. This network automatically adjusts depth of each stack with respect to the capability of its correspondent separable filter on extracting higher order convolutional features for the dataset under study. We conduct our experiments on several publicly available datasets varying in number of classes and quality of samples by using a standard implementation. Our results confirm the supremacy of our method on improving the precision of texture understanding in comparison with the recently published benchmarks.

## A Mathematics of Projection/Backprojection

We start by formulating of a classical dimension reduction problem and extend its solution to our proposed supervised projection/backprojection paradigms.

Given a sample set  $\mathcal{X} = \{X_1, X_2, \dots, X_{|X|}\}$  in  $\mathbb{R}^d$ , we try to find a matrix  $\mathbf{B} \in \mathbb{R}^{d \times c}$  that maps the input vector  $x_i$  onto the point  $y_i = \mathbf{B}^T x_i$  in a lower dimensional space  $\mathbb{R}^c$  conditioned on  $c \ll d$  by maximizing the separability between and minimizing the scattering within classes of set  $\mathcal{X}$ .

One of the solutions is supervised learning via linear discriminant analysis (LDA) [2]. The mapping matrix  $\mathbf{B}$  is determined to maximize the Fisher criterion given by

$$\mathcal{J}_{\mathcal{F}}(\mathbf{B}) = \text{tr} \left( (\mathbf{B} \mathbf{S}_{\mathbf{wB}} \mathbf{B}^T)^{-1} (\mathbf{B} \mathbf{S}_{\mathbf{bB}} \mathbf{B}^T) \right) \quad (1)$$

which  $tr(\cdot)$  is diagonal summation operator. The within/between-class scatterings  $\{\mathbf{S}_{\mathbf{w}B}, \mathbf{S}_{\mathbf{b}B}\}$  are defined as

$$\mathbf{S}_{\mathbf{w}B} = \sum_{j=1}^c \sum_{x_i \in \mathbf{C}_j} (x_i - \mu_j)(x_i - \mu_j)^T \quad (2)$$

$$\mathbf{S}_{\mathbf{b}B} = \sum_{j=1}^c (\mu_j - \bar{\mu})(\mu_j - \bar{\mu})^T \quad (3)$$

where  $c$ ,  $\mu_j$  and  $\bar{\mu}$  are number of classes, mean over class  $\mathbf{C}_j$  and mean over all dataset, respectively.

The matrix  $\mathbf{S}_{\mathbf{w}B}$  can be regarded as the average class-specific covariance, whereas  $\mathbf{S}_{\mathbf{b}B}$  can be viewed as the mean distance between all different classes. Thus, the purpose of Equation 1 is to maximize the between-class scatter while preserving within-class dispersion.

The answer is the solution of generalized eigenvalue problem  $\mathbf{S}_{\mathbf{b}B} \mathbf{B} = \lambda \mathbf{S}_{\mathbf{w}B} \mathbf{B}$ . Since the rank of  $\mathbf{S}_{\mathbf{b}B}$  is  $c - 1$ , the solution for  $c$  classes is eigenvectors corresponding to the largest  $c - 1$  eigenvalues of  $\mathbf{S}_{\mathbf{w}B}^{-1} \mathbf{S}_{\mathbf{b}B}$  for  $c \ll d$  [10].

Assuming  $\mathbf{S}_{\mathbf{w}B} \neq \mathbf{I}$  and  $\mathbf{S}_{\mathbf{b}B} \neq \mathbf{I}$ , by cyclic permutation of trace operator and imposing orthogonality through  $\mathbf{B} \mathbf{B}^T = \mathbf{I}$  (identity matrix), the Equation 1 holds

$$\begin{aligned} \mathcal{J}_{\mathcal{F}}(\mathbf{B}) &= tr\left(\left(\mathbf{B} \mathbf{S}_{\mathbf{w}B} \mathbf{B}^T\right)^{-1} \left(\mathbf{B} \mathbf{S}_{\mathbf{b}B} \mathbf{B}^T\right)\right) \\ &= tr\left(\left(\mathbf{B}^T\right)^{-1} \mathbf{S}_{\mathbf{w}B}^{-1} \mathbf{B}^{-1} \mathbf{B} \mathbf{S}_{\mathbf{b}B} \mathbf{B}^T\right) \\ &= tr\left(\mathbf{B}^T \left(\mathbf{B}^T\right)^{-1} \mathbf{S}_{\mathbf{w}B}^{-1} \times \mathbf{I} \times \mathbf{S}_{\mathbf{b}B}\right) \\ &= tr\left(\mathbf{S}_{\mathbf{w}B}^{-1} \mathbf{S}_{\mathbf{b}B}\right) \end{aligned} \quad (4)$$

To come up with our proposed projection for  $c > d$ , we again consider the Fisher criterion in Equation 1 and redefine inter-class scattering  $\mathbf{S}_{\mathbf{w}A} \in \mathbb{R}^{c \times c}$  such that it satisfies

$$tr(\mathbf{S}_{\mathbf{w}B}) = tr(\mathbf{S}_{\mathbf{w}A}) \quad (5)$$

Note that in Equations 2, we sum over all classes ( $c$ ) and hence, to satisfy Equation 5, we can consider  $\mathbf{S}_{\mathbf{w}A}$  as a square matrix of size  $c \times c$  with all zeros except main diagonal entries

$$\mathbf{S}_{\mathbf{w}A}(j, j) = tr\left(\sum_{x_i \in \mathbf{C}_j} (x_i - \mu_j)(x_i - \mu_j)^T\right) \quad \forall j \in [1, c] \quad (6)$$

From Equations 5 and similarity invariance of trace operator,  $\mathbf{S}_{\mathbf{w}B}$  and  $\mathbf{S}_{\mathbf{w}A}$  are similar matrices [11] which implies, there should exist a non-singular matrix  $\mathbf{\Gamma}_{\mathbf{w}}$  such that

$$\mathbf{S}_{\mathbf{w}B} = \mathbf{\Gamma}_{\mathbf{w}}^{-1} \mathbf{S}_{\mathbf{w}A} \mathbf{\Gamma}_{\mathbf{w}} \quad (7)$$

By minor matrix operations, Equation 7 can be formed as

$$\mathbf{\Gamma}_{\mathbf{w}} \mathbf{S}_{\mathbf{w}B} - \mathbf{S}_{\mathbf{w}A} \mathbf{\Gamma}_{\mathbf{w}} = \mathbf{0} \quad (8)$$

which is a special case of Sylvester equation [4] and can be solved for  $\mathbf{\Gamma}_w$  by either Kronecker tensor trick or using generalized eigen decomposition because, we define  $\mathbf{S}_{wB}$  and  $\mathbf{S}_{wA}$  as non-singular matrices. The closed form solution for Equation 8 is

$$\text{vec}(\mathbf{\Gamma}_w) = \mathbf{I} \otimes (-\mathbf{S}_{wA}) - \mathbf{S}_{wB}^T \otimes \mathbf{I} \quad (9)$$

which  $\text{vec}(\cdot)$  is vectorization operator. With the same reasoning, we define  $\mathbf{S}_{bA}$  as a square matrix of size  $c \times c$  such that

$$\text{tr}(\mathbf{S}_{bB}) = \text{tr}(\mathbf{S}_{bA}) \quad (10)$$

and there should exist a non-singular matrix  $\mathbf{\Gamma}_b$  such that

$$\mathbf{S}_{bB} = \mathbf{\Gamma}_b^{-1} \mathbf{S}_{bA} \mathbf{\Gamma}_b \quad (11)$$

On the other hand, from Equations 7 and 11

$$\begin{aligned} \mathbf{S}_{wB}^{-1} \mathbf{S}_{bB} &= (\mathbf{\Gamma}_w^{-1} \mathbf{S}_{wA} \mathbf{\Gamma}_w)^{-1} (\mathbf{\Gamma}_b^{-1} \mathbf{S}_{bA} \mathbf{\Gamma}_b) \\ &= \mathbf{\Gamma}_w^{-1} \mathbf{S}_{wA}^{-1} \mathbf{\Gamma}_w \mathbf{\Gamma}_b^{-1} \mathbf{S}_{bA} \mathbf{\Gamma}_b \end{aligned} \quad (12)$$

Due to the similarity invariance in Equations 5 and 10, we consider the cyclic permutation of trace operator and suppose that

$$\mathbf{\Gamma}_b = \mathbf{\Gamma}_w \quad (13)$$

and hence, Equation 11 implies  $\mathbf{S}_{bA}$  as

$$\mathbf{S}_{bA} = \mathbf{\Gamma}_b \mathbf{S}_{bB} \mathbf{\Gamma}_b^{-1} \quad (14)$$

Now, we work out Equation 12 by substitution from Equation 13 as follows

$$\begin{aligned} \mathbf{S}_{wB}^{-1} \mathbf{S}_{bB} &= \mathbf{\Gamma}_w^{-1} \mathbf{S}_{wA}^{-1} \times \mathbf{I} \times \mathbf{S}_{bA} \mathbf{\Gamma}_b \\ &= \mathbf{\Gamma}_w^{-1} (\mathbf{S}_{wA}^{-1} \mathbf{S}_{bA}) \mathbf{\Gamma}_w \end{aligned} \quad (15)$$

that proves  $\mathbf{S}_{wB}^{-1} \mathbf{S}_{bB}$  and  $\mathbf{S}_{wA}^{-1} \mathbf{S}_{bA}$  are similar matrices such that it holds

$$\text{tr}(\mathbf{S}_{wB}^{-1} \mathbf{S}_{bB}) = \text{tr}(\mathbf{S}_{wA}^{-1} \mathbf{S}_{bA}) \quad (16)$$

Looking back at Equation 4, we are able to define a new optimization problem for  $\mathbf{S}_{wA}$  and  $\mathbf{S}_{bA}$  considering the same discrimination power and orthogonality of Equation 1 as

$$\mathcal{J}_{\mathcal{F}}(\mathbf{A}) = \text{tr} \left( (\mathbf{A} \mathbf{S}_{wA} \mathbf{A}^T)^{-1} (\mathbf{A} \mathbf{S}_{bA} \mathbf{A}^T) \right) \quad (17)$$

which is aligned with the number of classes ( $c$ ) instead of the dimension of the input ( $d$ ). Employing the same eigenvector solution of Equation 1 to maximize Equation 17, gives the projection matrix  $\mathbf{A} \in \mathbb{R}^{d \times c}$  for  $c > d$ .

## B Closed Forms of Gradients

Suppose that  $\mathcal{H}(\mathbf{A})$  is composed of  $\mathcal{H}_1(\mathbf{A})$  and  $\mathcal{H}_2(\mathbf{A})$  as follows

$$\mathcal{H}_1(\mathbf{A}) = \frac{\text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{wA}} \mathbf{A}^T)}{\text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{bA}} \mathbf{A}^T)} \quad (18)$$

$$\mathcal{H}_2(\mathbf{A}) = \|\mathbf{I} - \mathbf{A} \mathbf{A}^T\|_2 \quad (19)$$

According to matrix calculus [9],

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{wA}} \mathbf{A}^T)}{\partial \mathbf{A}} = (\mathbf{S}_{\mathbf{wA}}^T + \mathbf{S}_{\mathbf{wA}}) \mathbf{A}^T \quad (20)$$

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{bA}} \mathbf{A}^T)}{\partial \mathbf{A}} = (\mathbf{S}_{\mathbf{bA}}^T + \mathbf{S}_{\mathbf{bA}}) \mathbf{A}^T \quad (21)$$

and hence, we have

$$\begin{aligned} \frac{\partial \mathcal{H}_1}{\partial \mathbf{A}} &= \frac{(\mathbf{S}_{\mathbf{bA}}^T + \mathbf{S}_{\mathbf{bA}}) \mathbf{A}^T \times \text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{wA}} \mathbf{A}^T)}{\left(\text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{bA}} \mathbf{A}^T)\right)^2} \\ &\quad - \frac{(\mathbf{S}_{\mathbf{wA}}^T + \mathbf{S}_{\mathbf{wA}}) \mathbf{A}^T \times \text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{bA}} \mathbf{A}^T)}{\left(\text{tr}(\mathbf{A} \mathbf{S}_{\mathbf{bA}} \mathbf{A}^T)\right)^2} \end{aligned} \quad (22)$$

On the other hand,

$$\frac{\partial \mathcal{H}_2}{\partial \mathbf{A}} = \frac{\partial(\mathbf{I} - \mathbf{A} \mathbf{A}^T)}{\partial \mathbf{A}} \times \frac{\mathbf{I} - \mathbf{A} \mathbf{A}^T}{\|\mathbf{I} - \mathbf{A} \mathbf{A}^T\|_2} \quad (23)$$

which gives

$$\frac{\partial \mathcal{H}_2}{\partial \mathbf{A}} = \frac{-2 \mathbf{A}^T \times (\mathbf{I} - \mathbf{A} \mathbf{A}^T)}{\|\mathbf{I} - \mathbf{A} \mathbf{A}^T\|_2} \quad (24)$$

The derivatives for  $\mathcal{Q}(\mathbf{B})$  can be calculated with the same reasoning as

$$\begin{aligned} \frac{\partial \mathcal{Q}_1}{\partial \mathbf{B}} &= \frac{(\mathbf{S}_{\mathbf{bB}} + \mathbf{S}_{\mathbf{bB}}^T) \mathbf{B} \times \text{tr}(\mathbf{B}^T \mathbf{S}_{\mathbf{wB}} \mathbf{B})}{\left(\text{tr}(\mathbf{B}^T \mathbf{S}_{\mathbf{bB}} \mathbf{B})\right)^2} \\ &\quad - \frac{(\mathbf{S}_{\mathbf{wB}} + \mathbf{S}_{\mathbf{wB}}^T) \mathbf{B} \times \text{tr}(\mathbf{B}^T \mathbf{S}_{\mathbf{bB}} \mathbf{B})}{\left(\text{tr}(\mathbf{B}^T \mathbf{S}_{\mathbf{bB}} \mathbf{B})\right)^2} \end{aligned} \quad (25)$$

$$\frac{\partial \mathcal{Q}_2}{\partial \mathbf{B}} = \frac{-2 \mathbf{B} \times (\mathbf{I} - \mathbf{B}^T \mathbf{B})}{\|\mathbf{I} - \mathbf{B}^T \mathbf{B}\|_2} \quad (26)$$

## References

- [1] Christopher M Bishop. Pattern recognition. *Machine Learning*, 2006.
- [2] Reinosuke FUKUNAGA. Statistical pattern recognition. 1990.
- [3] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- [4] Sang-Gu Lee and Quoc-Phong Vu. Simultaneous solutions of sylvester equations and idempotent matrices separating the joint spectrum. *Linear Algebra and its Applications*, 435(9):2097–2109, 2011.
- [5] Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7:15, 2008.