

# Stochastic Dykstra Algorithms for Metric Learning on Positive Semi-Definite Cone

Tomoki Matsuzawa<sup>†</sup>, Raissa Relator<sup>◊</sup>, Jun Sese<sup>◊</sup>, Tsuyoshi Kato<sup>†,‡,\*</sup>

<sup>†</sup> Faculty of Science and Engineering, Gunma University, Kiryu-shi, Gunma, 326-0338, Japan.

<sup>‡</sup> Center for Informational Biology, Ochanomizu University, Bunkyo-ku, Tokyo, 112-8610, Japan. n

<sup>◊</sup> BRD, AIST, Koto-ku, Tokyo, 135-0064, Japan.

## Abstract

Recently, covariance descriptors have received much attention as powerful representations of set of points. In this research, we present a new metric learning algorithm for covariance descriptors based on the Dykstra algorithm, in which the current solution is projected onto a half-space at each iteration, and runs at  $O(n^3)$  time. We empirically demonstrate that randomizing the order of half-spaces in our Dykstra-based algorithm significantly accelerates the convergence to the optimal solution. Furthermore, we show that our approach yields promising experimental results on pattern recognition tasks.

## 1 Introduction

Learning with example objects characterized by a set of several points, instead of a single point, in a feature space is an important task in the computer vision and pattern recognition community. In the case of visual categorization of still images, many local image descriptors such as SIFT [17] are extracted from an input image to form a single vector such as a Bag-of-Visual-Words vector or a Fisher Vector [18, 20]. For image set classification, a surge of methods have been developed in the last decade, and probabilistic models [25] or kernels [22] are introduced to describe the image set. Alternative descriptors are the covariance descriptors, which have received much attention as a powerful representation of a set of points.

The performance of categorizing covariance descriptors depends on the metric that is used to measure the distances between them. To compare covariance descriptors, a variety of distance measures such as affine invariant Riemannian metric [19], Stein metric [27], J-divergence [29], Frobenius distance [11], and Log-Frobenius distance [1], have been discussed in existing literature. Some of them are designed from their geometrical properties, but some are not. Many of these distance measures are expressed in the form

$$D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2) := \|\Phi(\mathbf{X}_1) - \Phi(\mathbf{X}_2)\|_F^2,$$

with  $\Phi : \mathbb{S}_{++}^n \rightarrow \mathbb{R}^{n \times m}$  for some  $m \in \mathbb{N}$ . If  $\Phi(\mathbf{X}) := \text{logm}(\mathbf{X})$ , where  $\text{logm}(\mathbf{X})$  takes the principal matrix logarithm of a strictly positive definite matrix  $\mathbf{X}$ , the Log-Frobenius distance [1] is obtained. Setting  $\Phi(\mathbf{X}) := \mathbf{X}^p$  gives the Power-Frobenius distance [11], while  $\Phi(\mathbf{X}) := \text{chol}(\mathbf{X})$ , where  $\text{chol} : \mathbb{S}_{++}^n \rightarrow \mathbb{R}^{n \times n}$  produces the Cholesky decomposition of  $\mathbf{X}$  such that  $\mathbf{X} = \text{chol}(\mathbf{X})\text{chol}(\mathbf{X})^\top$ , yields the Cholesky-Frobenius distance [8]. These metrics are pre-defined before the employment of machine learning algorithms, and are not adaptive to the data to be analyzed. Meanwhile, for categorization of vectorial data, supervised learning for fitting metrics to the task has been proven to significantly increase the performance of the distance-based classifier [6, 13, 23].

In this paper, we introduce a parametric distance measure between covariance descriptors and present novel metric learning algorithms to determine the parameters of the distance measure function. The learning problem is formulated as the Bregman projection onto the intersections of half-spaces. This kind of problem can be solved by the Dykstra algorithm [4, 9], which chooses a single half-space in a cyclic order and projects a current solution to the half-space. We developed an efficient technique for projection onto a single half-space. Furthermore, we empirically found that selecting the half-space stochastically, rather than in a cyclic order, dramatically increases the speed of converging to an optimal solution.

### 1.1 Related work

To the best of our knowledge, Vemulapalli et al. (2015) [28] were the first to introduce the supervised metric learning approach for covariance descriptors. They vectorized the matrix logarithms of the covariance descriptors to apply existing metric learning methods to the vectorizations of matrices. The dimensionality of the vectorizations is  $n(n+1)/2$  when the size of the covariance matrices are  $n \times n$ . Thus, the size of the Mahalanobis matrix is  $n(n+1)/2 \times n(n+1)/2$ , which is computationally prohibitive when  $n$  is large.

Our approach is an extension of the distance measure of Huang et al. [10], which is based on the Log-Euclidean metric, with their loss function being a special case of our formulation. They also adopted the cyclic Dykstra algorithm for learning the Mahalanobis-like matrix. However, they misused the Woodbury matrix inversion formula when deriving the projection onto a single half-space, therefore, their algorithm has no theoretical guarantee of converging to the optimal solution. In this paper, their update rule is corrected by presenting a new technique that projects a current solution to a single half-space within  $O(n^3)$  computational time.

Yger and Sugiyama [30] devised a different formulation of metric learning. They introduced the congruent transform [2] and measures distances between the transformations of covariance descriptors. An objective function based on the kernel target alignment [5] is employed to determine the transformation parameters. Compared to their algorithm, our algorithm has the capability to monitor the upper bound of the objective gap, i.e. the difference between the current objective and the minimum. This implies that the resultant solution is ensured to be  $\epsilon$ -suboptimal if the algorithm's convergence criterion is set such that the objective gap upper bound is less than a very small number  $\epsilon$ . Since Yger and Sugiyama [30] employed a gradient method for learning the congruent transform, there is no way to know the objective gap.

## 1.2 Contributions

Our contributions of this paper can be summarized as follows.

- For metric learning on positive semidefinite cone, we developed a new algorithm based on the Dykstra algorithm, in which the current solution is projected onto a half-space at each iterate, and runs at  $O(n^3)$  time.
- We present an upper-bound for the objective gap which provides a stopping criterion and ensures the optimality of the solution.
- We empirically found that randomizing the order of half-spaces in our Dykstra-based algorithm significantly accelerates the convergence to the optimal solution.
- We show that our approach yields promising experimental results on pattern recognition tasks.

## 1.3 Notation

We denote vectors by bold-faced lower-case letters and matrices by bold-faced upper-case letters. Entries of vectors and matrices are not bold-faced. The transposition of a matrix  $\mathbf{A}$  is denoted by  $\mathbf{A}^\top$ , and the inverse of  $\mathbf{A}$  is by  $\mathbf{A}^{-1}$ . The  $n \times n$  identity matrix is denoted by  $\mathbf{I}_n$ . The subscript is often omitted. The  $m \times n$  zero matrix is denoted by  $\mathbf{O}_{m \times n}$ . The subscript is often omitted. The  $n$ -dimensional vector all of whose entries are one is denoted by  $\mathbf{1}_n$ . We use  $\mathbb{R}$  and  $\mathbb{N}$  to denote the set of real and natural numbers,  $\mathbb{R}^n$  and  $\mathbb{N}^n$  to denote the set of  $n$ -dimensional real and natural vectors, and  $\mathbb{R}^{m \times n}$  to denote the set of  $m \times n$  real matrices. For any  $n \in \mathbb{N}$ , we use  $\mathbb{N}_n$  to denote the set of natural numbers less than or equal to  $n$ . Let us define  $\mathbb{R}_+ := \{x \in \mathbb{R} \mid x \geq 0\}$ ,  $\mathbb{R}_{++} := \{x \in \mathbb{R} \mid x > 0\}$ ,  $\mathbb{R}_+^n := \{x \in \mathbb{R}^n \mid \mathbf{x} \geq \mathbf{0}_p\}$ , and  $\mathbb{R}_{++}^n := \{x \in \mathbb{R}^n \mid \mathbf{x} > \mathbf{0}_p\}$ . The relational operator  $\succ$  denotes the generalized inequality associated with the strictly positive definite cone. We use  $\mathbb{S}^n$  to denote the set of symmetric  $n \times n$  matrices.  $\mathbb{S}_+^n$  to denote the set of symmetric positive semi-definite  $n \times n$  matrices, and  $\mathbb{S}_{++}^n$  to denote the set of symmetric strictly positive definite  $n \times n$  matrices. For any  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$ ,  $\text{diag}(\mathbf{x})$  is defined as an  $n \times n$  diagonal matrix whose diagonal entries are  $x_1, \dots, x_n$ . For any  $n \times n$  square matrix  $\mathbf{X}$ , its trace is denoted by  $\text{tr}(\mathbf{X})$ . For any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , define  $\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{i=1}^n x_i y_i$  where  $x_i$  and  $y_i$  is the  $i$ -th entry of  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. For any  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}$ , define  $\langle \mathbf{X}, \mathbf{Y} \rangle := \sum_{i=1}^m \sum_{j=1}^n X_{i,j} Y_{i,j}$  where  $X_{i,j}$  and  $Y_{i,j}$  is the  $(i, j)$ -th entry of  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.  $\mathbb{O}_n$  is used to denote the set of  $n \times n$  orthonormal matrices, i.e.  $\mathbb{O}_n := \{\mathbf{A} \in \mathbb{R}^{n \times n} \mid \mathbf{A}^\top \mathbf{A} = \mathbf{I}_n\}$ .

# 2 Our Metric Learning Problem

## 2.1 Parametric distance measure on $\mathbb{S}_+^n$

We introduce the following distance measure for covariance descriptors  $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{S}_+^n$ :

$$D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2; \mathbf{W}) := \left\langle \mathbf{W}, (\Phi(\mathbf{X}_1) - \Phi(\mathbf{X}_2)) (\Phi(\mathbf{X}_1) - \Phi(\mathbf{X}_2))^\top \right\rangle,$$

where  $\mathbf{W} \in \mathbb{S}_+^n$  is the parameter of this distance measure function. If  $\mathbf{W}$  is strictly positive definite and  $\Phi$  is bijective, then this distance measure  $D_{\Phi}(\cdot, \cdot; \mathbf{W}) : \mathbb{S}_+^n \times \mathbb{S}_+^n \rightarrow \mathbb{R}$  is a metric because all of the following conditions are satisfied: (i) non-negativity:  $D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2; \mathbf{W}) \geq 0$ ; (ii) identity of indiscernibles:  $D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2; \mathbf{W}) = 0$  iff  $\mathbf{X}_1 = \mathbf{X}_2$ ; (iii) symmetry:  $D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2; \mathbf{W}) = D_{\Phi}(\mathbf{X}_2, \mathbf{X}_1; \mathbf{W})$ ; (iv) triangle inequality:  $D_{\Phi}(\mathbf{X}_1, \mathbf{X}_3; \mathbf{W}) \leq D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2; \mathbf{W}) + D_{\Phi}(\mathbf{X}_2, \mathbf{X}_3; \mathbf{W})$ .

If the parameter matrix  $\mathbf{W}$  is singular,  $D_{\Phi}(\cdot, \cdot; \mathbf{W})$  is a pseudometric, and the identity of indiscernibles is changed to the following property: For any  $\mathbf{X}_1 \in \mathbb{S}_{++}^n$ ,  $D_{\Phi}(\mathbf{X}_1, \mathbf{X}_1; \mathbf{W}) = 0$  holds, while  $D_{\Phi}(\mathbf{X}_1, \mathbf{X}_2; \mathbf{W}) = 0$  occurs for some non-identical positive semi-definite matrices  $\mathbf{X}_1$  and  $\mathbf{X}_2$ .

## 2.2 Formulations of the learning problems

To determine the value of the parameter matrix  $\mathbf{W}$ , we pose a constrained optimization problem based on the idea of ITML [6]. We now consider a multi-class categorization problem. Let  $n_c$  be the number of classes, and the class labels are represented by natural numbers in  $\mathbb{N}_{n_c}$ . Suppose we are given

$$(\mathbf{X}_1, \omega_1), \dots, (\mathbf{X}_\ell, \omega_\ell) \in \mathbb{S}_+^n \times \mathbb{N}_{n_c}$$

as a training dataset, where  $\mathbf{X}_i$  is the covariance descriptor of the  $i$ -th example, and  $\omega_i$  is its class label. From the  $\ell$  examples,  $K$  pairs  $(i_1, j_1), \dots, (i_K, j_K) \in \mathbb{N}_\ell \times \mathbb{N}_\ell$  are picked to give, to each pair, the following constraint:

$$D_{\Phi}(\mathbf{X}_{i_k}, \mathbf{X}_{j_k}; \mathbf{W}) \begin{cases} \leq b_{\text{ub}} \xi_k, & \text{if } \omega_{i_k} = \omega_{j_k}, \\ \geq b_{\text{lb}} \xi_k, & \text{if } \omega_{i_k} \neq \omega_{j_k}, \end{cases} \quad (1)$$

where, when  $\xi_k = 1$ , the two constants  $b_{\text{ub}}$  and  $b_{\text{lb}}$ , respectively, are the upper-bound of the distances between any two examples in the same class and the lower-bound of the distances between any two examples in different classes. Now let us define for  $k \in \mathbb{N}_K$ ,

$$y_k := \begin{cases} +1, & \text{if } \omega_{i_k} = \omega_{j_k}, \\ -1, & \text{if } \omega_{i_k} \neq \omega_{j_k}, \end{cases} \quad \text{and} \quad b_k := \begin{cases} b_{\text{ub}}, & \text{if } \omega_{i_k} = \omega_{j_k}, \\ b_{\text{lb}}, & \text{if } \omega_{i_k} \neq \omega_{j_k}. \end{cases}$$

Under the constraint (1), we wish to find  $\mathbf{W}$  and  $\xi_k$  such that  $\mathbf{W}$  is not much deviated from the identity matrix and  $\xi_k$  is close to one. From this motivation, we pose the following problem:

$$\begin{aligned} \min \quad & \text{BD}_{\varphi}((\mathbf{W}, \boldsymbol{\xi}), (\mathbf{I}, \mathbf{1})), \quad \text{wrt } \mathbf{W} \in \mathbb{S}_{++}^n, \quad \boldsymbol{\xi} = [\xi_1, \dots, \xi_K]^{\top} \in \mathbb{R}_{++}^K, \\ \text{subject to} \quad & \forall k \in \mathbb{N}_K, \quad y_k D_{\Phi}(\mathbf{X}_{i_k}, \mathbf{X}_{j_k}; \mathbf{W}) \leq y_k b_k \xi_k, \end{aligned} \quad (2)$$

where  $\text{BD}_{\varphi}(\cdot, \cdot) : (\mathbb{S}_{++}^n \times \mathbb{R}_{++}^K) \times (\mathbb{S}_{++}^n \times \mathbb{R}_{++}^K) \rightarrow \mathbb{R}_+$  is the *Bregman divergence* [14]. Only if  $(\mathbf{W}, \boldsymbol{\xi}) = (\mathbf{I}, \mathbf{1})$  will the divergence  $\text{BD}_{\varphi}((\mathbf{W}, \boldsymbol{\xi}), (\mathbf{I}, \mathbf{1}))$  become zero, and the value of divergence becomes larger if  $(\mathbf{W}, \boldsymbol{\xi})$  is more deviated from  $(\mathbf{I}, \mathbf{1})$ . The definition of the Bregman divergence contains a seed function  $\varphi : \mathbb{S}_{++}^n \times \mathbb{R}_{++}^K \rightarrow \mathbb{R}$  which is assumed to be continuously differentiable and strictly convex. For some  $\varphi$ , the Bregman divergence is defined as

$$\text{BD}_{\varphi}(\boldsymbol{\Theta}, \boldsymbol{\Theta}_0) = \varphi(\boldsymbol{\Theta}) - \varphi(\boldsymbol{\Theta}_0) - \langle \nabla \varphi(\boldsymbol{\Theta}_0), \boldsymbol{\Theta} - \boldsymbol{\Theta}_0 \rangle,$$

for  $\boldsymbol{\Theta}, \boldsymbol{\Theta}_0 \in \mathbb{S}_{++}^n \times \mathbb{R}_{++}^K$ . This implies that the quantities of the deviations of the solution  $(\mathbf{W}, \boldsymbol{\xi})$  from  $(\mathbf{I}, \mathbf{1})$  depend on the definition of the seed function. In this study, the seed function is assumed to be the sum of two terms:

$$\varphi(\mathbf{W}, \boldsymbol{\xi}) := \varphi_{\text{r}}(\mathbf{W}) + \sum_{k=1}^K c_k \varphi_1(\xi_k),$$

where  $c_k$  is a positive constant. The first term  $\varphi_{\text{r}} : \mathbb{S}_{++}^n \rightarrow \mathbb{R}$  in the definition of the seed function is defined by  $\varphi_{\text{r}}(\mathbf{W}) := -\log \det(\mathbf{W})$ . As for the definition of the second term  $\varphi_1 : \mathbb{R}_{++} \rightarrow \mathbb{R}$ , we considered the following three functions:

$$\varphi_{\text{is}}(\xi_k) := -\log(\xi_k), \quad \varphi_{\text{l2}}(\xi_k) := \frac{1}{2} \xi_k^2, \quad \varphi_{\text{e}}(\xi_k) := (\log \xi_k - 1) \xi_k.$$

The Bregman divergences generated from three seed functions  $\varphi_{\text{is}}$ ,  $\varphi_{\text{l2}}$ , and  $\varphi_{\text{e}}$ , respectively, are referred to as *Itakura-Saito Bregman Divergence* (ISBD), *L2 Bregman Divergence* (L2BD), and *Relative Entropy Bregman Divergence* (REBD), where ISBD is equal to the objective function employed by Huang et al. [10].

## 3 Stochastic Variants of Dykstra Algorithm

We introduce the Dykstra algorithm [4, 9] to solve the optimization problem (2). The original Dykstra algorithm [9] was developed as a computational method that finds the Euclidean projection from a point onto the intersection of

convex sets. Censor & Reich [4] extended the algorithm to finding the Bregman projection from a point  $\mathbf{x}_0$  to a set  $\mathcal{C}$ , defined by

$$\operatorname{argmin}_{\mathbf{x} \in \mathcal{C}} \operatorname{BD}_\varphi(\mathbf{x}, \mathbf{x}_0).$$

In available literature related to stochastic gradient descent methods and the variants [3, 12, 24, 26] that minimize the regularized loss averaged over a set of examples, it is empirically shown that, rather than picking an example in a cyclic order, example selection in a stochastic order dramatically speeds up the convergence to the optimal solution. Alternatively, some literature reported that at the beginning of every epoch in the gradient method, random permutation of the order of examples also accelerates the convergence [7].

Motivated by these facts, this study proposes the use of stochastic orders for selection of convex set components in the Dykstra algorithm. We term the stochastic version of the Dykstra algorithm as the *stochastic Dykstra algorithm*. In our case, every convex set component is one of  $K$  half-spaces, as will be described in a later discussion. There are, then, three ways to select half-spaces:

- **Cyclic:** Pick a half-space in a cyclic order at each iteration.
- **Rand:** Pick a half-space randomly at each iteration.
- **Perm:** Permute the order of  $K$  half-spaces randomly at the beginning of each epoch.

Hereinafter, we assume to employ the “Rand” option, although replacing this option with one of the remaining two is straightforward.

If every convex set component is a half-space, and the  $k$ -th convex set component  $\mathcal{C}_k$  is expressed as

$$\mathcal{C}_k := \{\mathbf{x} \mid \langle \mathbf{a}_k, \mathbf{x} \rangle \leq b_k\},$$

then computing the Bregman projection from a point  $\mathbf{x}_0$  to its boundary  $\operatorname{bd}(\mathcal{C}_k)$  is equivalent to solving the following saddle point problem:

$$\max_{\delta} \min_{\mathbf{x}} \operatorname{BD}_\varphi(\mathbf{x}, \mathbf{x}_0) + \delta(\langle \mathbf{a}_k, \mathbf{x} \rangle - b_k).$$

This fact enables us to rewrite the Dykstra algorithm with Rand option for finding the Bregman projection from a point  $\mathbf{x}_0$  to the intersection of  $\mathcal{C}_1, \dots, \mathcal{C}_K$ , as described in Algorithm 1, where  $\varphi^*$  is the convex conjugate of the seed function  $\varphi$ .

---

**Algorithm 1** Stochastic Dykstra Algorithm.

---

- 1: **begin**
- 2:  $\forall k \in \mathbb{N}_K : \alpha_k := 0;$
- 3: **for**  $t = 1, 2, \dots$  **do**
- 4:   Pick  $k$  randomly from  $\{1, \dots, K\};$
- 5:   Solve the following saddle point problem and let  $\delta_{t-1/2}$  be the solution of  $\delta:$

$$\max_{\delta} \min_{\mathbf{x}} \operatorname{BD}_\varphi(\mathbf{x}, \mathbf{x}_{t-1}) + \delta_t(\langle \mathbf{a}_k, \mathbf{x} \rangle - b_k); \tag{3}$$

- 6:    $\delta_t := \max(\delta_{t-1/2}, -\alpha_k); \alpha_k := \alpha_k + \delta_t;$
  - 7:    $\mathbf{x}_t = \nabla \varphi^*(\nabla \varphi(\mathbf{x}_{t-1}) - \delta_t \mathbf{a}_k);$
  - 8: **end for**
  - 9: **end.**
- 

## 4 Efficient Projection Technique

We now show that solving the optimization problem (2) is equivalent to finding a Bregman projection from a point  $(\mathbf{I}, \mathbf{1}) \in \mathbb{S}_{++}^n \times \mathbb{R}_{++}^K$  onto the intersection of multiple half-spaces.

Let  $\mathbf{A}_k$  be a positive semidefinite matrix expressed as

$$\mathbf{A}_k := (\Phi(\mathbf{X}_{i_k}) - \Phi(\mathbf{X}_{j_k})) (\Phi(\mathbf{X}_{i_k}) - \Phi(\mathbf{X}_{j_k}))^\top$$

for  $k \in \mathbb{N}_K$ , to define a half-space

$$\mathcal{C}_k := \{(\mathbf{W}, \boldsymbol{\xi}) \in \mathbb{S}_{++}^n \times \mathbb{R}_{++}^K \mid y_k \langle \mathbf{A}_k, \mathbf{W} \rangle - y_k b_k \xi_k \leq 0\}.$$

Then, it can be seen that the intersection of  $K$  half-spaces

$$\bigcap_{k=1}^K \mathcal{C}_k$$

is the feasible region of the optimization problem (2). This implies that the Dykstra algorithm can be applied to solve problem (2).

Next we present an efficient technique that projects  $(\mathbf{W}_{t-1}, \boldsymbol{\xi}_{t-1}) \in \mathbb{S}_{++}^n \times \mathbb{R}_{++}^K$  onto the  $k$ -th half-space  $\mathcal{C}_k$ , where  $(\mathbf{W}_{t-1}, \boldsymbol{\xi}_{t-1}) \in \mathbb{S}_{++}^n \times \mathbb{R}_{++}^K$  is the model parameter after the  $(t-1)$ -th iteration. Let  $\xi_{k,t-1}$  be the  $k$ -th entry in the vector  $\boldsymbol{\xi}_{t-1}$ . The value of the function  $J_t : \mathbb{R} \rightarrow \mathbb{R}$  defined by

$$J_t(\delta) := \langle \mathbf{A}_k, (\mathbf{W}_{t-1}^{-1} + \delta y_k \mathbf{A}_k)^{-1} \rangle - b_k \nabla \varphi_1^*(\nabla \varphi_1(\boldsymbol{\xi}_{t-1}) + \delta y_k b_k / c_k),$$

is zero at the solution  $\delta_{t-1/2}$  of the saddle point problem (3). The solution  $\delta$  must satisfy the strictly positive definiteness:

$$(\mathbf{Y}(\delta))^{-1} := \mathbf{W}_{t-1}^{-1} + \delta y_k \mathbf{A}_k \succ \mathbf{O}, \quad (4)$$

and the feasibility of the slack variables:

$$\exists \xi_{k,t-1/2} \quad \text{s.t.} \quad \nabla \varphi_1(\xi_{k,t-1/2}) = \nabla \varphi_1(\xi_{k,t-1}) - \delta y_k b_k / c_k. \quad (5)$$

There is no closed-form solution found for this projection problem. Hence, some numerical method such as the Newton-Raphson method is necessary for solving the nonlinear system  $J_t(\delta) = 0$ . If one tries to compute the value of  $J_t(\delta)$  naïvely, it will require an  $O(n^3)$  computational cost because  $J_t(\cdot)$  involves computation of the inverse of an  $n \times n$  matrix. If we suppose the numerical method assesses the value of the scalar-valued function  $J_t(\cdot)$   $L$  times, the naïve approach will take  $O(Ln^3)$  computational time to find the solution of the nonlinear system  $J_t(\delta) = 0$ . Furthermore, the positive definiteness condition in (4) and the feasibility condition in (5) must be checked.

We will show the following two claims:

- The solution of the system  $J_t(\delta) = 0$  satisfying (4) and (5) can be computed within  $O(n^3 + Ln)$  time, where  $L$  is the number of times a numerical method assesses the value of  $J_t(\cdot)$ .
- The solution exists and it is unique.

Hereinafter, we assume  $\mathbf{A}_k$  is strictly positive definite. By setting  $\mathbf{A}_k \leftarrow \mathbf{A}_k + \epsilon \mathbf{I}$ , with  $\epsilon$  as a small positive constant, it is easy to satisfy this assumption. Since  $L \in O(n^2)$  in a typical setting, we can say that each update can be done in  $O(n^3)$  computation.

We define  $\mathbf{A}_k^{1/2}$ ,  $\mathbf{U}$ ,  $\mathbf{D}$ , and  $\mathbf{d}$  as follows. Let  $\mathbf{A}_k^{1/2} \in \mathbb{S}_{++}^n$  such that  $\mathbf{A}_k^{1/2} \mathbf{A}_k^{1/2} = \mathbf{A}_k$ , and denote by  $\mathbf{A}_k^{-1/2} \in \mathbb{S}_{++}^n$  the inverse of  $\mathbf{A}_k^{1/2}$ . Introduce an orthonormal matrix  $\mathbf{U} \in \mathbb{O}_n$  and a diagonal matrix  $\mathbf{D} = \text{diag}(\{d_1, \dots, d_n\})$  that represent a spectral decomposition  $\mathbf{U} \mathbf{D} \mathbf{U}^\top = \mathbf{A}_k^{-1/2} \mathbf{W}_{t-1}^{-1} \mathbf{A}_k^{-1/2}$ , with  $d_1 \geq \dots \geq d_n$ . Then, we have

$$\mathbf{Y}(\delta) = \mathbf{A}_k^{-1/2} \mathbf{U} (\mathbf{D} + y_k \delta \mathbf{I})^{-1} \mathbf{U}^\top \mathbf{A}_k^{-1/2}, \quad (6)$$

which allows us to rewrite the first term of  $J_t(\delta)$  as

$$\langle \mathbf{A}_k, (\mathbf{W}_{t-1}^{-1} + \delta y_k \mathbf{A}_k)^{-1} \rangle = \sum_{i=1}^n \frac{1}{d_i + y_k \delta}. \quad (7)$$

Assessment of  $J_t(\delta)$  can be done within  $O(n)$  computational cost after  $d_1, \dots, d_n$  are obtained. To get the  $n$  scalars  $d_1, \dots, d_n$ , we need to find  $\mathbf{A}_k^{-1/2}$  and the spectral decomposition of  $\mathbf{A}_k^{-1/2} \mathbf{W}_{t-1}^{-1} \mathbf{A}_k^{-1/2}$ , each of which requires  $O(n^3)$  computation. The  $n \times n$  matrix  $\mathbf{A}_k^{-1/2}$  can be computed in the pre-process of the Dykstra algorithm, while the spectral decomposition of  $\mathbf{A}_k^{-1/2} \mathbf{W}_{t-1}^{-1} \mathbf{A}_k^{-1/2}$  is done once before invoking some numerical method to solve the nonlinear system  $J_t(\delta) = 0$ . These support the first claim.

Equation (6) suggests that the set of  $\delta$  satisfying (4) is given by

$$I_{r,t} := \begin{cases} (-d_n, +\infty), & \text{for } y_k = +1, \\ (-\infty, d_n), & \text{for } y_k = -1. \end{cases} \quad (8)$$

The set of  $\delta$  satisfying (5) is given as follows. In the case of using ISBD,  $\delta$  ensuring (5) is in the interval

$$I_{\text{is},t} := \begin{cases} (-\infty, \delta_{\text{b}}), & \text{for } y_k = +1, \\ (-\delta_{\text{b}}, +\infty), & \text{for } y_k = -1, \end{cases}$$

where

$$\delta_{\text{b}} := \frac{c_k}{b_k \xi_{k,t-1}}.$$

In the case of using L2BD and REBD, there exists  $\xi_{t-1/2}$  even if  $\nabla\varphi_1(\xi_{t-1}) - \delta y_k b_k / c_k$  takes any value.

Hence, if ISBD is employed, the solution  $\delta_{t-1/2}$  can be searched from the interval

$$I_{\text{r},t} \cap I_{\text{is},t} = \begin{cases} (-d_n, \delta_{\text{b}}), & \text{for } y_k = +1, \\ (-\delta_{\text{b}}, +d_n), & \text{for } y_k = -1. \end{cases}$$

If L2BD or REBD is employed, the solution  $\delta_{t-1/2}$  can be searched from  $I_{\text{r},t}$ . In the reminder of this section, we shall use the notation  $I_t$  to denote the interval for  $\delta$  satisfying (4) and (5) simultaneously.

We now show the uniqueness of the solution. The gradient of  $J_t : \mathbb{R} \rightarrow \mathbb{R}$  is expressed as

$$\nabla J_t(\delta) = - \sum_{i=1}^n \frac{y_k}{(d_i + y_k \delta)^2} - \frac{y_k b_k^2}{c_k} \nabla^2 \varphi_1^* (\nabla \varphi_1(\xi_{t-1}) - \delta y_k b_k / c_k),$$

for  $\delta \in I_t$ . We first consider the case that  $y_k = +1$ . Clearly, the first term is negative. The second term is non-positive because any convex conjugate function is convex. Therefore, we have  $\nabla J_t(\delta) < 0$ . In the case of  $y_k = -1$ , we get  $\nabla J_t(\delta) > 0$  from a similar derivation. These observations imply that the solution is unique if a solution exists. The existence of the solution can be established by showing that the curve  $J_t(\delta)$  crosses the horizontal axis.

We consider the cases of using ISBD and using either L2BD or REBD separately. For the ISBD case, we have

$$\lim_{\delta \searrow -d_n} J_t(\delta) = +\infty, \quad \lim_{\delta \nearrow \delta_{\text{b}}} J_t(\delta) = -\infty,$$

if  $y_k = +1$ , and

$$\lim_{\delta \searrow -\delta_{\text{b}}} J_t(\delta) = -\infty, \quad \lim_{\delta \nearrow d_n} J_t(\delta) = +\infty,$$

if  $y_k = -1$ . On the other hand, when using either L2BD or REBD with  $y_k = +1$  we get

$$\lim_{\delta \searrow -d_n} J_t(\delta) = +\infty, \quad \lim_{\delta \rightarrow +\infty} J_t(\delta) = -\infty,$$

while we obtain

$$\lim_{\delta \rightarrow -\infty} J_t(\delta) = -\infty, \quad \lim_{\delta \nearrow d_n} J_t(\delta) = +\infty,$$

when  $y_k = -1$ . Hence, we conclude that

$$\exists! \delta \in I_t \quad \text{s.t.} \quad J_t(\delta) = 0.$$

## 4.1 Stopping Criterion

Here we discuss how to determine if the solution is already optimal and when to terminate the algorithm. While running the algorithm,  $(\mathbf{W}_t, \boldsymbol{\xi}_t)$  may be infeasible to the primal problem. Denote the index set of the violated constraints by  $\mathcal{I}_{\text{vio}} := \{k \in \mathbb{N}_K \mid (\mathbf{W}_t, \boldsymbol{\xi}_t) \notin \mathcal{C}_k\}$  and let us define  $\bar{\boldsymbol{\xi}}_t \in \mathbb{R}_{++}^K$  so that the  $k$ -th entry is given by  $\bar{\xi}_{h,t} := \frac{1}{b_h} \langle \mathbf{W}_t, \mathbf{A}_h \rangle$  for  $h \in \mathcal{I}_{\text{vio}}$  and  $\bar{\xi}_{h,t} := \xi_{h,t}$  for  $h \notin \mathcal{I}_{\text{vio}}$ . Note that  $(\mathbf{W}_t, \bar{\boldsymbol{\xi}}_t)$  is a feasible solution, and  $\bar{\boldsymbol{\xi}}_t = \boldsymbol{\xi}_t$  when  $(\mathbf{W}_t, \boldsymbol{\xi}_t)$  is feasible. The objective gap after iteration  $t$  is bounded as follows:

$$\text{BD}_{\varphi}((\mathbf{W}_t, \bar{\boldsymbol{\xi}}_t), (\mathbf{I}, \mathbf{1})) - \text{BD}_{\star} \leq \sum_{h \in \mathcal{I}_{\text{vio}}} c_h (\varphi_1(\bar{\xi}_{h,t}) - \varphi_1(\xi_{h,t}) - \nabla \varphi_1(1)(\bar{\xi}_{h,t} - \xi_{h,t})) - \sum_{h=1}^K \alpha_h y_h (\langle \mathbf{A}_h, \mathbf{W}_t \rangle - b_h \xi_{h,t}),$$

where we have defined

$$\text{BD}_{\star} := \min_{(\mathbf{W}, \boldsymbol{\xi}) \in \bigcap_h \mathcal{C}_h} \text{BD}_{\varphi}((\mathbf{W}, \boldsymbol{\xi}), (\mathbf{I}, \mathbf{1})).$$

Then this upper-bound of the objective gap can be used for the stopping criterion of the Dykstra algorithm.

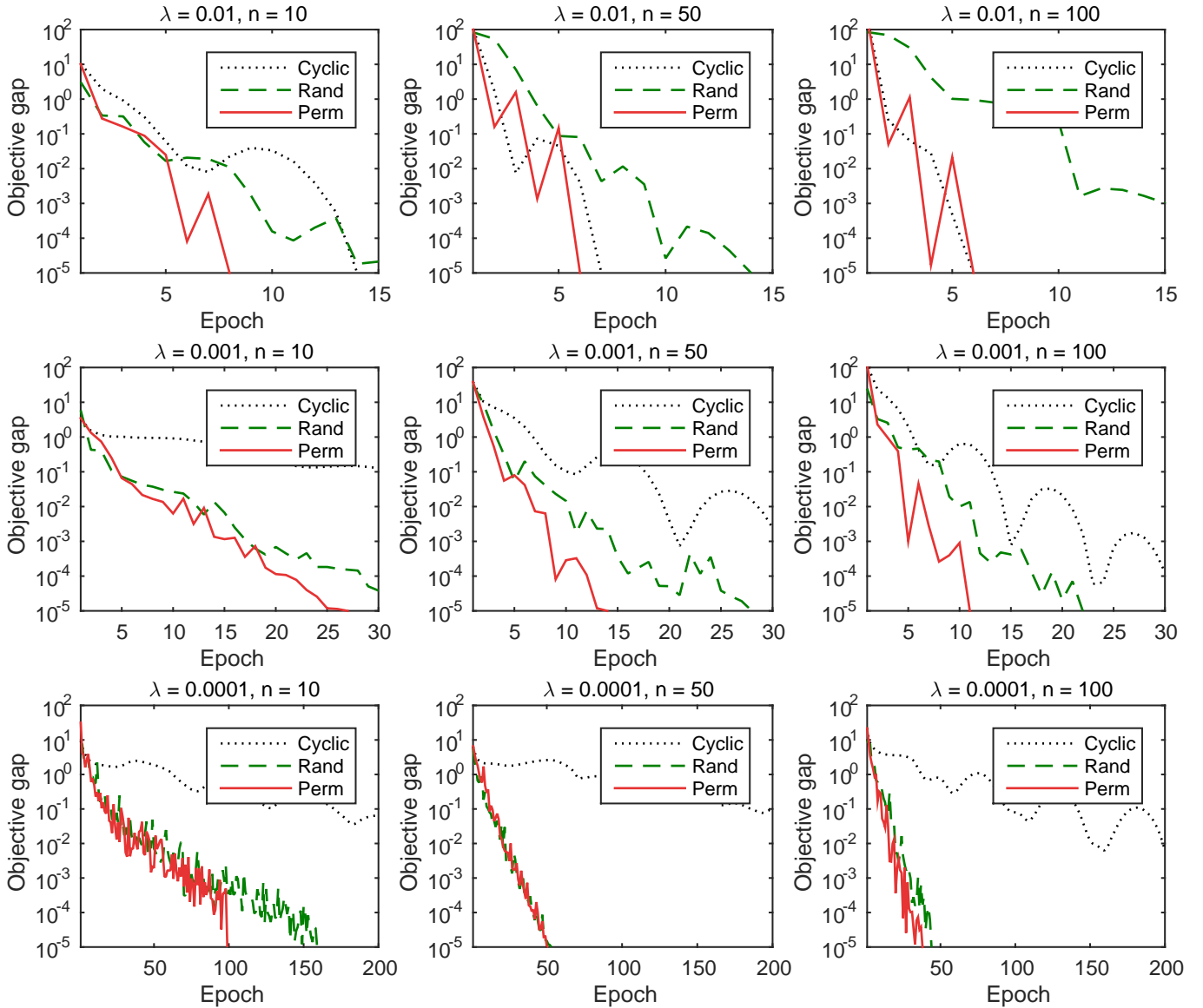


Figure 1: Convergence behavior of the algorithms using different settings.

## 5 Numerical Experiments

We conducted experiments to assess the convergence speed of our optimization algorithms and the generalization performance for pattern recognition.

### 5.1 Convergence behavior of optimization algorithms

We examined our algorithms for assessment of convergence speed. We generated datasets artificially as follows.  $K = 50$  matrices  $\mathbf{F}_k \in \mathbb{R}^{n \times n}$  are generated in which each entry is drawn from the uniform distribution in the interval  $[-0.5, 0.5]$ . Then, we set  $\mathbf{A}_k := \mathbf{F}_k \mathbf{F}_k^\top$ . The values of the variables  $y_k$  are randomly chosen from  $\{\pm 1\}$  with same probabilities. We set  $\mathbf{b} = \mathbf{1}$  and  $\mathbf{c} = \mathbf{1}/(\lambda K)$ . We exhaustively tested Cyclic, Perm, and Rand with the settings of  $\lambda = 10^{-2}, 10^{-3}, 10^{-4}$  and  $n = 10, 50, 100$ .

Figure 1 demonstrates the convergence behavior of the cyclic Dykstra algorithm and the two stochastic Dykstra algorithm with various  $\lambda$  and  $n$ . Here, one epoch is called  $K$  times projection onto a single half-space. ISBD is employed as the objective function for learning the metric  $\mathbf{W}$ . The objective gap is defined as the difference between the current objective value and the minimum. In most of the settings, the two stochastic Dykstra algorithms converged faster than the cyclic algorithm. Especially when  $\lambda = 10^{-4}$ , the cyclic algorithm was too slow to use it in practice.

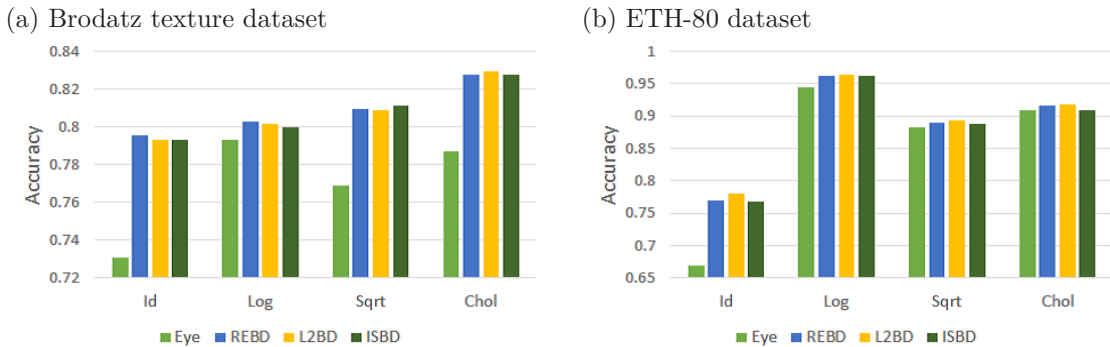


Figure 2: Generalized performances for pattern recognition.

## 5.2 Generalization performance for pattern recognition

We used the Brodatz texture dataset [21] containing 111 different texture images to examine the generalization performance for texture classification. Each image has a size of  $640 \times 640$  and gray-scaled. Images were individually divided into four sub-images of equal size. One of the four sub-images was picked randomly and used for testing, and the rest of the images were used for training.

For each training image and each testing image, covariance descriptors of randomly chosen 50 were extracted from  $128 \times 128$  patches. The covariance matrices are of five-dimensional feature vectors  $[I, |I_x|, |I_y|, |I_{xx}|, |I_{yy}|]^T$ . Then, 11,100 ( $= 111 \times 2 \times 50$ ) covariance descriptors are obtained for training and testing, respectively. For evaluation of generalized performance,  $k$ -nearest neighbor classifier is used, where the number of the nearest neighbors is set to three. We set  $K = 100 \times n_c$ ,  $b_{ub} = 0.05$ , and  $b_{lb} = 0.95$ .

We also examined the generalization performance for generic visual categorization using the ETH-80 dataset [16] containing  $n_c = 8$  classes. Each class has 10 objects, each of which includes 41 colored images. For every object, 20 images are randomly chosen and used for training, and the rest of images are used for testing.

One covariance matrix is obtained from each image. Eight features  $[x, y, R, G, B, |I_x|, |I_y|, |I_{xx}|, |I_{yy}|]^T$  are obtained from each pixel in an image.

We tried four types of  $\Phi$ : **Id**:  $\Phi(\mathbf{X}) = \mathbf{X}$ , **Log**:  $\Phi(\mathbf{X}) = \log_m(\mathbf{X})$ , **Sqrt**:  $\Phi(\mathbf{X}) = \mathbf{X}^{1/2}$ , **Chol**:  $\Phi(\mathbf{X}) = \text{chol}(\mathbf{X})$ . The parameter  $\mathbf{W}$  is determined by the metric learning algorithms with ISBD, L2BD, and REBD, to be compared with  $\mathbf{W} = \mathbf{I}$  we denote as **Eye**. Note that  $D_{\Phi}(\cdot, \cdot; \mathbf{I}) = D_{\Phi}(\cdot, \cdot)$ . Figure 2 gives the accuracy bar plots for the two multi-class classification problems. Whichever  $\Phi$  is used, supervised metric learning improved the generalization performances both for texture classification and for generic visual categorization. For texture classification, the Cholesky decomposition-based mapping  $\text{chol}(\cdot)$  achieved the best accuracy, while the matrix logarithm-based mapping  $\log_m(\cdot)$  obtained the highest accuracy for generic image categorization.

## 6 Conclusions

In this paper, we have devised several objective functions for metric learning on positive semidefinite cone, all of which can be minimized by the Dykstra algorithm. We have introduced a new technique that performs each update efficiently when the Dykstra algorithm is applied to the metric learning problems. We have empirically demonstrated that the stochastic versions of the Dykstra algorithm are much faster than the original algorithm.

## Acknowledgment

This work was supported by JSPS KAKENHI Grant Number 26249075, 40401236. The last author would like to thank Dr. Zhiwu Huang for fruitful discussions.

## References

- [1] Vincent Arsigny, Pierre Fillard, Xavier Pennec, and Nicholas Ayache. Log-euclidean metrics for fast and simple calculus on diffusion tensors. *Magnetic resonance in medicine*, 56(2):411–421, 2006.
- [2] Rajendra Bhatia. *Positive Definite Matrices*. Princeton University Press, 2009.



- [3] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In Yves Lechevallier and Gilbert Saporta, editors, *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT'2010)*, pages 177–187, Paris, France, August 2010. Springer.
- [4] Yair Censor and Simeon Reich. The dykstra algorithm with bregman projections. *Communications in Applied Analysis*, 2:407–419, 1998.
- [5] Nello Cristianini, John Shawe-Taylor, André Elisseeff, and Jaz S. Kandola. On kernel-target alignment. In Thomas G. Dietterich, Suzanna Becker, and Zoubin Ghahramani, editors, *NIPS*, pages 367–373. MIT Press, 2001.
- [6] Jason V. Davis, Brian Kulis, Prateek Jain, Suvrit Sra, and Inderjit S. Dhillon. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, pages 209–216. ACM, 2007.
- [7] Aaron J Defazio, Tibério S Caetano, and Justin Domke. Finito: A faster, permutable incremental gradient method for big data problems. *arXiv preprint arXiv:1407.2710*, 2014.
- [8] Ian L Dryden, Alexey Koloydenko, and Diwei Zhou. Non-euclidean statistics for covariance matrices, with applications to diffusion tensor imaging. *The Annals of Applied Statistics*, pages 1102–1123, 2009.
- [9] Richard L. Dykstra. An algorithm for restricted least squares regression. *Journal of the American Statistical Association*, 78(384):837–842, December 1983.
- [10] Zhiwu Huang, Ruiping Wang, Shiguang Shan, Xianqiu Li, and Xilin Chen. Log-euclidean metric learning on symmetric positive definite manifold with application to image set classification. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015*, pages 720–729, 2015.
- [11] Sadeep Jayasumana, Richard Hartley, Mathieu Salzmann, Hongdong Li, and Mehrtaf Harandi. Kernel methods on the riemannian manifold of symmetric positive definite matrices. In *CVPR*, pages 73–80. IEEE, 2013.
- [12] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems 26: Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 315–323, 2013.
- [13] Tsuyoshi Kato and Nozomi Nagano. Metric learning for enzyme active-site search. *Bioinformatics*, 26(21):2698–2704, November 2010.
- [14] Tsuyoshi Kato, Wataru Takei, and Shinichiro Omachi. A discriminative metric learning algorithm for face recognition. *IPSJ Transactions on Computer Vision and Applications*, 5:85–89, 2013.
- [15] Tsuyoshi Kato, Koji Tsuda, and Kiyoshi Asai. Selective integration of multiple biological data for supervised network inference. *Bioinformatics*, 21:2488–2495, May 2005.
- [16] Bastian Leibe and Bernt Schiele. Analyzing appearance and contour based methods for object categorization. In *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, volume 2, pages II–409. IEEE, 2003.
- [17] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.
- [18] Tomoki Matsuzawa, Raissa Relator, Wataru Takei, Shinichiro Omachi, and Tsuyoshi Kato. Mahalanobis encodings for visual categorization. *IPSJ Transactions on Computer Vision and Applications*, 7:69–73, July 2015. doi: 10.2197/ipsjtcva.7.1.
- [19] Xavier Pennec, Pierre Fillard, and Nicholas Ayache. A riemannian framework for tensor computing. *International Journal of Computer Vision*, 66(1):41–66, 2006.
- [20] Florent Perronnin, Jorge Sanchez, and Thomas Mensink. Improving the fisher kernel for large-scale image classification. In *Computer Vision–ECCV 2010*, pages 143–156. Springer, 2010.
- [21] Trygve Randen and John Hakon Husoy. Filtering for texture classification: A comparative study. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 21(4):291–310, 1999.
- [22] Raissa Relator, Yoshihiro Hirohashi, Eisuke Ito, and Tsuyoshi Kato. Mean polynomial kernel and its application to vector sequence recognition. *IEICE Transactions on Information and Systems*, E97-D(7):1855–1863, July 2014.

- [23] Raissa Relator, Nozomi Nagano, and Tsuyoshi Kato. Using bregmann divergence regularized machine for comparison of molecular local structures. *IEICE Transactions on Information & Systems*, E99-D(1):-, Jan 2016.
- [24] Nicolas L. Roux, Mark Schmidt, and Francis R. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 2663–2671. Curran Associates, Inc., 2012.
- [25] Gregory Shakhnarovich, John W Fisher, and Trevor Darrell. Face recognition from long-term observations. In *ECCV 2002*, pages 851–865. Springer Berlin Heidelberg, 2002.
- [26] Shai Shalev-Shwartz, Yoram Singer, Nathan Srebro, and Andrew Cotter. Pegasos: primal estimated sub-gradient solver for SVM. *Math. Program.*, 127(1):3–30, 2011.
- [27] Suvrit Sra. A new metric on the manifold of kernel matrices with application to matrix geometric means. In *Advances in Neural Information Processing Systems*, pages 144–152, 2012.
- [28] Raviteja Vemulapalli and David W Jacobs. Riemannian metric learning for symmetric positive definite matrices. *arXiv preprint arXiv:1501.02393*, 2015.
- [29] Zhizhou Wang and Baba C Vemuri. An affine invariant tensor dissimilarity measure and its applications to tensor-valued image segmentation. In *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, volume 1, pages I–228. IEEE, 2004.
- [30] Florian Yger and Masashi Sugiyama. Supervised logeuclidean metric learning for symmetric positive definite matrices. *arXiv preprint arXiv:1502.03505*, 2015.