

Positive Definite Matrices: Data Representation and Applications to Computer Vision

Anoop Cherian and Suvrit Sra

Abstract Numerous applications in computer vision and machine learning rely on representations of data that are compact, discriminative, and robust while satisfying several desirable invariances. One such recently successful representation is offered by *symmetric positive definite* (SPD) matrices. However, the modeling power of SPD matrices comes at a price: rather than a flat Euclidean view, SPD matrices are more naturally viewed through curved geometry (Riemannian or otherwise) which often complicates matters. We focus on models and algorithms that rely on the geometry of SPD matrices, and make our discussion concrete by casting it in terms of covariance descriptors for images. We summarize various commonly used distance metrics on SPD matrices, before highlighting formulations and algorithms for solving sparse coding and dictionary learning problems involving SPD data. Through empirical results, we showcase the benefits of mathematical models that exploit the curved geometry of SPD data across a diverse set of computer vision applications.

1 Introduction

Efficient representations that compactly capture salient properties of data form the basis of every algorithm in computer vision and machine learning. Consider for instance the task of tracking a person across video frames given (i) a video sequence, and (ii) a bounding box around the object to be tracked (see Figure 1). Methods for this task typically first generate a representation (descriptor) of the image patch inside the bounding box and then proceed further. Success of the tracking application relies on several desirable properties, for instance (i) ability of the representation to

Anoop Cherian
ARC Centre of Excellence for Robotic Vision, Australian National University, Canberra, Australia
<http://www.roboticvision.org> e-mail: anoop.cherian@anu.edu.au

Suvrit Sra
Massachusetts Institute of Technology, Cambridge, MA, USA e-mail: suvrit@mit.edu

uniquely characterize the contents of the tracked patch against any other patch in a video frame; (ii) robustness of the descriptor to camera sensor noise; (iii) stability despite changes in illuminations; (iv) tolerance for occlusions; and (v) robustness to affine deformations of the object.

A simple descriptor for tracking could be the normalized color histogram of an image patch. This choice partially satisfies properties such as (i), (ii), and (iii). We can improve this naive color descriptor by adding additional statistics of the image patch, such as the shape and pose of the object, texture of the patch, edge orientations, etc. However, each of these additional data features requires high-dimensional descriptors for its representation, whereby the final augmented data descriptor can be extremely large. This in turn raises storage concerns, while also complicating learning, recognition, and tracking due to the curse of dimensionality.



Fig. 1 An illustrative application: people tracking in a video sequence. We are given a video sequence and a bounding box of the person to be tracked (dark box), and the problem is to track the person along the sequence, i.e., to generate the lighter (yellow) boxes in subsequent frames.

A simple but effective alternative representation that fuses such multi-modal cues was introduced in [37], dubbed *region covariance descriptors*. The key idea is to model an image patch using correlations between different low-level features. To this end, we extract raw image features such as color, intensity gradients, etc., from every pixel within the desired patch. Then we stack these raw features and compute their covariance matrix, resulting in the *covariance descriptor* (for the patch).

Although it seems surprising that such a simple approach can lead to a powerful descriptor, we must note that the diagonal of the covariance matrix captures the statistical variance of each individual feature and the off-diagonals capture the correlation between different features. Thus, this descriptor captures second-order co-occurrences of multimodal features. Since this matrix has size quadratic only in the number of features, it is independent of the number of pixels in the patch and is thus very compact. Also, since the features means are subtracted when computing the covariance, this descriptor implicitly applies mean-filtering to the data, providing noise robustness. Due to these unique properties, the use of covariance descriptors has proliferated into several applications in computer vision and beyond.

We note that SPD matrices also play an important role as data descriptors in several other non-vision applications, such as, diffusion tensor imaging [3], brain-computer interfaces [17], sound compression [45], polarimetric image modeling [24], virus classification [21], and as quantum density matrices [20]. While, some of these matrices are not covariances, they are indeed symmetric positive definite and the algorithms proposed in this chapter are directly applicable.

While constructing covariance descriptors is simple (see Section 1.1), using them can be demanding. The key difficulty arises due to the non-Euclidean nature of covariances. Although it is tempting to simply view SPD matrices using the geometry of their embedding (Euclidean) space [2], the corresponding manifold is not compact which can lead to irrelevant solutions for certain applications [4]. Therefore, several alternative geometries for SPD matrices have been considered—Section 1.2 outlines some of the more widely used choices.

Subsequently, in Section 2 we illustrate use of these geometries on the concrete tasks of sparse coding and dictionary learning. These problems have been widely studied (for vectors) and are important to a variety of applications [34]. However, their extension to SPD matrix data is non-trivial and less studied; we describe two different frameworks that achieve these extensions.

1.1 Covariance Descriptors and Example Applications

Before delving into theoretical details, we recall below details on construction of covariance descriptors and summarize two illustrative applications.

Definition 1 (Covariance Descriptor). Suppose we have a data instance $\mathbf{I} \in \mathbb{R}^n$ (image patch, video snippet, etc.). Let $f_1^j, f_2^j, \dots, f_d^j$, (each $f_i^j \in \mathbb{R}, j = 1, 2, \dots, n$) represent d features computed for the j -th component of \mathbf{I} (such as image intensity gradients, filter outputs, etc.). Let $\mathbf{f}^j = [f_1^j, f_2^j, \dots, f_d^j]^T$; then the *Covariance Descriptor* S for \mathbf{I} is the $d \times d$ covariance matrix given by:

$$S = \frac{1}{n} \sum_{j=1}^n (\mathbf{f}^j - \boldsymbol{\mu})(\mathbf{f}^j - \boldsymbol{\mu})^T, \quad (1)$$

where $\boldsymbol{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{f}^j$ is the mean feature vector. We will assume that the features are linearly independent, in which case $S \in \mathbb{S}_+^d$, the cone of $d \times d$ SPD matrices.

1.1.1 Illustrative Application 1: People Tracking

We continue with our example of people tracking described in Figure 1. One important choice that we must make is what features to use for constructing the descriptor. For people tracking, we will first consider a texture based approach as described in [37] that uses raw image intensity features from each pixel. Given a patch \mathbf{I} , let $I_r(x, y), I_g(x, y), I_b(x, y)$ represent the red, green, and blue color intensities respec-

tively of pixel at spatial coordinates (x, y) in the patch. Further, let $I_x(x, y), I_y(x, y)$ represent the gray scale image intensity gradients in the horizontal and vertical directions respectively. Then, we define

$$\mathbf{f}(x, y) = [I_r(x, y), I_g(x, y), I_b(x, y), I_x(x, y), I_y(x, y)]^T \quad (2)$$

as the feature vector. A covariance descriptor capturing color and shape of the objects in \mathbf{I} is then given by:

$$S_{tracking} = \frac{1}{|\mathbf{I}|} \sum_{x, y} (\mathbf{f}(x, y) - \boldsymbol{\mu})(\mathbf{f}(x, y) - \boldsymbol{\mu})^T, \quad (3)$$

where $\boldsymbol{\mu} = \frac{1}{|\mathbf{I}|} \sum_{x, y} \mathbf{f}(x, y)$ is the mean feature, and $|\mathbf{I}|$ is the patch size. Note that the order of the pixels in the patch are ignored when computing the covariance matrix, thus providing invariance to changes in the pose of the object. However, we may associate additional spatial correlations between the pixel features by including the (x, y) coordinates as well into (2). In case, rotational invariance is desired, instead of including (x, y) coordinates, we could include $\sqrt{x^2 + y^2}$ as an additional feature for a pixel at location (x, y) . Curvature information could be easily included by computing the second-order gradients for every pixel and augmenting (2). This illustration shows the flexibility of the covariance descriptor to blend any feature of choice into the framework easily. The next example shows a richer set of features.

1.1.2 Illustrative Application 2: Face Recognition

While using raw features such as gradients could be sufficient to capture textures, recognizing faces require much more expressive and discriminative feature sets. In [39], the authors propose to use covariance descriptors for this task where the features are generated from Gabor wavelets. These filters are believed to reflect the human visual system [35] and measures the energy distribution in the image patch at various scales and orientations. These filters have been previously shown to be useful to extract discriminative and subtle features suitable for recognition [31]; using them within a covariance descriptor setup is shown to lead to significantly better recognition performance in [39], via capturing their second-order statistics. To this end, 40 Gabor wavelet filters are first designed consisting of 8 orientations and 5 different scales. Next, small patches centered at each pixel in the face patch are convolved with these filters, thus forming a 40 dimensional feature vector for each pixel, which precedes applying the steps described above to generate 40×40 covariance descriptors.

Covariance descriptors have been found to be useful in several other vision applications, such as visual surveillance [33, 36, 51], object recognition [22, 27], action recognition [47, 25, 50], image set classification [52], and emotion classification [53], to name a few. Almost all these works use a similar setup for computing the covariances, except that they use various features suitable for the application.

1.2 Geometry of SPD Matrices

To effectively use covariance descriptors we must next define schemes to compute similarity / distance between two descriptors. Intuitively, the distance measure enforces some geometry on the space of these descriptors, and which particular geometry is preferable depends on the application. As alluded to previously, a variety of SPD geometries have been explored, the simplest of which is the usual Euclidean geometry where the distance between SPD matrices \mathbf{X} and \mathbf{Y} is simply the Frobenius distance given by

$$d_F(\mathbf{X}, \mathbf{Y}) := \|\mathbf{X} - \mathbf{Y}\|_F. \quad (4)$$

However, d_F is neither affine invariant nor does it lead to a complete metric space (due to singular boundary). Affine invariance is important in applications such as diffusion MRI [40] while completeness is crucial when defining convergent sequences on the SPD manifold.

Two basic alternatives popular in computer vision are (i) the affine invariant Riemannian metric (AIRM) [40]; and (ii) the log-Euclidean Riemannian metric (LERM) [4]. Both measures induce a Riemannian geometry; the former induces a curved geometry while the latter “flattens” the manifold by mapping into the tangent space (which is Euclidean). The corresponding distance functions are

$$d_R(\mathbf{X}, \mathbf{Y}) := \left\| \text{Log} \mathbf{X}^{-1/2} \mathbf{Y} \mathbf{X}^{-1/2} \right\|_F, \quad (5)$$

$$d_{LE}(\mathbf{X}, \mathbf{Y}) := \|\text{Log} \mathbf{X} - \text{Log} \mathbf{Y}\|_F, \quad (6)$$

where $\mathbf{X}^{-1/2}$ is the matrix square-root of SPD matrix \mathbf{X}^{-1} and Log is the principal matrix logarithm. Distance (5) is affine invariant, and enjoys a host of other remarkable properties [7, Ch. 6]. LERM, however, is not affine invariant though it is rotation and scale invariant separately. Both AIRM and LERM are computationally demanding: for $d \times d$ SPD matrices, assuming approximately $4d^3/3$ flops for eigendecomposition and d^3 for matrix multiplication, these distances requires approximately $14d^3/3$ flops.

To reduce the computational cost, while preserving affine invariance and other geometric properties, the Stein distance was introduced in [15, 48]; it is defined as

$$d_S(\mathbf{X}, \mathbf{Y}) := \left[\log \det \left(\frac{\mathbf{X} + \mathbf{Y}}{2} \right) - \frac{1}{2} \log \det(\mathbf{X}\mathbf{Y}) \right]^{1/2}. \quad (7)$$

Computing d_S requires only three Cholesky decompositions, at a total cost of d^3 flops. Moreover, d_S is analytically simpler to work with, as its gradients are also simpler to obtain than either AIRM or LERM. Consequently, it has found application in a number of recent works, some of which we will refer to in the sequel.

2 Application to Sparse Coding and Dictionary Learning

After outlining a few basic SPD geometries, we are now ready to describe concrete applications where specific properties of SPD matrices play a role. In particular, we discuss sparse coding and dictionary learning for SPD valued data. The first model we cover is *generalized dictionary learning* (GDL), a direct extension of usual (vector) dictionary learning. Our second model uses the natural Riemannian geometry of SPD matrices, and measures sparse reconstruction loss using the AIRM distance. There are other formulations for dictionary learning and sparse coding on SPD matrices, for instance [25, 26, 46]; these differ from our model primarily in the formulation of the sparse reconstruction loss.

2.1 Dictionary Learning with SPD atoms

Traditional dictionary learning takes vectors $\mathbf{x}_i \in \mathbb{R}^p$ ($1 \leq i \leq m$) and constructs a matrix $\mathbf{B} \in \mathbb{R}^{p \times n}$ and code vectors $\boldsymbol{\alpha}_i \in \mathbb{R}^n$ (usually $n \gg p$), so that

$$\mathbf{x}_i \approx \mathbf{B}\boldsymbol{\alpha}_i, \quad \text{and } \boldsymbol{\alpha}_i \text{ is sparse, for } 1 \leq i \leq m.$$

The sparsity requirement on $\boldsymbol{\alpha}_i$ is commonly enforced using ℓ_0 - or ℓ_1 -norm penalties or constraints. Since both \mathbf{B} and $\boldsymbol{\alpha}_i$ are unknown dictionary learning usually results in a difficult nonconvex optimization task. Nevertheless, it has found remarkable success toward sparse coding and other applications [18].

We depart from the above setup in that we consider input matrices $\mathbf{X}_i \in \mathbb{R}^{p \times q}$, $1 \leq i \leq m$. Then, instead of a dictionary matrix \mathbf{B} we learn a tensor \mathbf{B} , which we identify with a linear operator $\mathbf{B} : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}^{p \times q}$ so that

$$\mathbf{X}_i \approx \mathbf{B}(\mathbf{A}_i), \quad \text{and } \mathbf{A}_i \text{ is sparse, for } 1 \leq i \leq m. \quad (8)$$

Using (8), one model of GDL is the following [49]:

$$\min_{\mathbf{A}_1, \dots, \mathbf{A}_m, \mathbf{B}} \frac{1}{2} \sum_{i=1}^m \|\mathbf{X}_i - \mathbf{B}(\mathbf{A}_i)\|_{\text{F}}^2 + \sum_{i=1}^m \beta_i \text{sp}(\mathbf{A}_i), \quad (9)$$

where $\beta_i > 0$ are scalar hyperparameters while $\text{sp}(\mathbf{A})$ enforces some notion of sparsity. For instance, $\text{sp}(\mathbf{A})$ could be the cardinality function $\|\mathbf{A}\|_0$, its convex relaxation $\|\mathbf{A}\|_1$, the matrix rank $\text{rank}(\mathbf{A})$ or its convex relaxation, the trace-norm $\|\mathbf{A}\|_{\text{tr}}$.

Formulation (9) requires one more modification for SPD valued inputs. Specifically, we ensure that the approximation $\mathbf{B}(\mathbf{A}_i)$ is also SPD by defining \mathbf{B} via

$$\mathbf{B}(\mathbf{A}) := \mathbf{B}\mathbf{A}\mathbf{B}^T, \quad \text{for some matrix } \mathbf{B}, \quad (10)$$

and additionally *restricting* to $\mathbf{A} \succeq 0$. Observe that (10) can be written as

$$\text{vec}(\mathbf{B}\mathbf{A}\mathbf{B}^T) = (\mathbf{B} \otimes \mathbf{B}) \text{vec}(\mathbf{A}), \quad (11)$$

where vec stacks columns of its argument and the operator \mathbf{B} is encoded (isomorphically) by the product $\mathbf{B} \otimes \mathbf{B}$.

It is easy to show that (11) requires $md^2 + d^2n + mn$ storage for m covariance matrices of size $d \times d$, while (10) takes $md^2 + dn + mn$. Computationally, also the first formulation is cheaper, so we prefer it. As to \mathbf{A} , we consider two choices:

1. $\mathbf{A} = \text{Diag}(\alpha_1, \dots, \alpha_n)$ where $\alpha_i \geq 0$; and
2. $\mathbf{A} = \sum_{j=1}^k \alpha_j \boldsymbol{\alpha}_j \boldsymbol{\alpha}_j^T$, a potentially low-rank (if $k < n$) SPD matrix.

Although diagonal \mathbf{A} might appear to be simple, it is quite powerful. Indeed, with it GDL models SPD matrices as weighted sums of rank-one matrices since

$$\mathbf{X} \approx \mathbf{B}\mathbf{A}\mathbf{B}^T = \sum_{i=1}^n \alpha_i \mathbf{b}_i \mathbf{b}_i^T, \quad \text{where } \alpha_i = \mathbf{A}_{ii}, \quad (12)$$

which offers a rich yet computationally tractable model.

2.1.1 Stochastic Gradient for GDL

Now we derive a stochastic-gradient procedure for approximately solving GDL. We use the convex function $\text{sp}(\mathbf{A}) = \|\mathbf{A}\|_1$ for enforcing sparsity. Then, using representation (12) with diagonal \mathbf{A}_i , the GDL optimization problem (9) becomes

$$\min_{\mathbf{A}_1, \dots, \mathbf{A}_m \geq 0, \mathbf{B}} \frac{1}{2} \sum_{i=1}^m \|\mathbf{X}_i - \mathbf{B}\mathbf{A}_i\mathbf{B}^T\|_{\mathbb{F}}^2 + \sum_{i=1}^m \beta_i \|\mathbf{A}_i\|_1. \quad (13)$$

Problem (13) is nonconvex and difficult. However, for a fixed dictionary \mathbf{B} , it is individually convex in $(\mathbf{A}_1, \dots, \mathbf{A}_m)$. It is thus amenable to the idea of alternating between updating \mathbf{B} and optimizing over $(\mathbf{A}_1, \dots, \mathbf{A}_m)$. But in many applications the number of input data points m is very large, so the alternating steps can easily become rather costly. Therefore, we follow a stochastic gradient approach that can scale to large data sets, as long as the stochastic gradients can be obtained efficiently.

To prevent degenerate solutions we also impose normalization constraints $\|\mathbf{b}_j\|_2 \leq 1$ on each column of matrix \mathbf{B} . We denote these requirements by the feasible set \mathcal{B} . We run stochastic gradient using K “mini-batches,” for which we rewrite (13) as

$$\min_{\mathbf{B} \in \mathcal{B}} \Phi(\mathbf{B}) := \sum_{b=1}^K \phi_b(\mathbf{B}), \quad (14)$$

where ϕ_b denotes the objective function for batch b . Let k_b be the size of batch b ($1 \leq b \leq K$) containing the matrices $\{\mathbf{X}_{j(i)} | 1 \leq i \leq k_b\}$, where $j(i)$ is an appropriate index in $1, \dots, m$. With this notation, the objective function for batch b is

$$\phi_b(\mathbf{B}) := \min_{\mathbf{A}_{j(1)}, \dots, \mathbf{A}_{j(k)} \geq 0} \frac{1}{2} \sum_{i=1}^{k_b} \|\mathbf{X}_{j(i)} - \mathbf{B}\mathbf{A}_{j(i)}\mathbf{B}^T\|_{\mathbb{F}}^2 + \beta_{j(i)} \|\mathbf{A}_{j(i)}\|_1. \quad (15)$$

We apply stochastic-gradient to (14), which performs the iteration

$$\mathbf{B}_{t+1} = \Pi_{\mathcal{B}}(\mathbf{B}_t - \eta_t \nabla_{\mathbf{B}} \phi_{b(t)}(\mathbf{B}_t)), \quad b(t) \in [1..K], \quad t = 0, 1, \dots, \quad (16)$$

where $\Pi_{\mathcal{B}}$ denotes orthogonal projection onto \mathcal{B} . Assuming (15) has a unique solution, the gradient $\nabla_{\mathbf{B}}\phi_{b(t)}$ is well defined. Specifically, let $(\mathbf{A}_{j(1)}^*, \dots, \mathbf{A}_{j(k)}^*)$ be the argmin of (15). Then, writing $b \equiv b(t)$, we have

$$\nabla_{\mathbf{B}}\phi_b(\mathbf{B}) = 2 \sum_{i=1}^{k_b} \left(\mathbf{B}\mathbf{A}_{j(i)}^* \mathbf{B}^T - \mathbf{X}_{j(i)} \right) \mathbf{B}\mathbf{A}_{j(i)}^*. \quad (17)$$

The computationally intensive part is to compute (17), which we now consider.

2.1.2 Sparse Coding: Computing $\nabla\phi_b$

Observe that (15) is a sum of k_b independent problems, so it suffices to describe the computation for a subproblem of the form

$$\min_{\mathbf{A} \geq 0} f(\mathbf{A}) := \frac{1}{2} \|\mathbf{X} - \mathbf{B}\mathbf{A}\mathbf{B}^T\|_{\mathbb{F}}^2 + \beta \|\mathbf{A}\|_1. \quad (18)$$

Since $\mathbf{A} \geq 0$ and diagonal, problem (18) is nothing but a regularized nonnegative least-squares (NNLS) problem. There exist a variety of solvers for NNLS, for example, LBFGS-B [32], or Spectral Projected-Gradient (SPG) [8]. We prefer to use the latter, as it is not only simple, but also exhibits excellent empirical performance.

In Section 3, we will apply this sparse coding scheme to the problem of nearest neighbor retrieval on covariance datasets. But, before proceeding to the experiments, we will elucidate a much richer and more powerful dictionary learning and sparse coding scheme on SPD matrices that leads to significantly better results on various applications; this scheme uses the natural Riemannian geometry for the sparse construction loss instead of the Euclidean distance as in (9).

2.2 Riemannian Dictionary Learning and Sparse Coding

Recall that we wish to compute a dictionary with ‘‘SPD atoms’’. We work on manifold $\mathbb{M}_n^d = \prod_{i=1}^n \mathbb{S}_+^d \subset \mathbb{R}^{d \times d \times n}$, which is the Cartesian product of n SPD manifolds. Our goals are (i) to learn a dictionary $\mathbf{B} \in \mathbb{M}_n^d$ in which each slice represents an SPD dictionary atom $\mathbf{B}_j \in \mathbb{S}_+^d$ $1 \leq j \leq n$; and (ii) to approximate each \mathbf{X}_i as a sparse conic combination of atoms in \mathbf{B} ; i.e., $\mathbf{X}_i \sim \mathbf{B}(\boldsymbol{\alpha}_i)$ where $\boldsymbol{\alpha}_i \in \mathbb{R}_+^n$ and $\mathbf{B}(\boldsymbol{\alpha}) := \sum_{i=1}^n \boldsymbol{\alpha}_i \mathbf{B}_i$. With this notation our *dictionary learning and sparse coding* (DLSC) problem is

$$\min_{\mathbf{B} \in \mathbb{M}_n^d, \boldsymbol{\alpha} \in \mathbb{R}_+^{n \times m}} \frac{1}{2} \sum_{j=1}^m \text{dR}^2(\mathbf{X}_j, \mathbf{B}\boldsymbol{\alpha}_j) + \text{Sp}(\boldsymbol{\alpha}_j) + \Omega(\mathbf{B}), \quad (19)$$

where Sp and Ω regularize the codes $\boldsymbol{\alpha}_j$ and the dictionary tensor \mathbf{B} respectively.

Formulation (19) is a direct SPD analog of the vector DL setup. Instead of learning a dictionary matrix for vectors, we learn a third-order tensor dictionary as our input is matricial data. We constraint the sparse codes to be non-negative to ensure that

the linear combination $\mathbf{B}(\boldsymbol{\alpha})$ remains SPD. In contrast to usual DL problems where the dictionary learning and sparse coding subproblems are convex, problem (19) is much harder: it is neither convex in itself nor are its subproblems convex.

Pragmatically speaking, this lack of subproblem convexity is not too damaging: we just need a set of dictionary atoms that can sparse code the input data, and such a set can still be computed by performing an alternating minimization (actually, just descent here). We describe the details below.

2.2.1 Dictionary Learning Subproblem

Assuming that the sparse code vectors $\boldsymbol{\alpha}$ are available, the subproblem of updating the dictionary atoms can be separated from (19) and written as:

$$\begin{aligned} \min_{\mathbf{B} \in \mathbb{M}_n^d} \Theta(\mathbf{B}) &:= \frac{1}{2} \sum_{j=1}^m d_{\mathbb{R}^2}^2(\mathbf{X}_j, \mathbf{B}\boldsymbol{\alpha}_j) + \Omega(\mathbf{B}), \\ &= \frac{1}{2} \sum_{j=1}^m \left\| \text{Log} \left(\mathbf{X}_j^{-1/2} (\mathbf{B}\boldsymbol{\alpha}_j) \mathbf{X}_j^{-1/2} \right) \right\|_F^2 + \Omega(\mathbf{B}). \end{aligned} \quad (20)$$

Due to its good empirical performance we choose $\Omega(\mathbf{B}) := \lambda_{\mathbf{B}} \sum_{i=1}^n \text{Tr}(\mathbf{B}_i)$.

Riemannian CG

Among several first-order alternatives for optimizing over the SPD atoms (such as the steepest-descent, trust-region methods [1], etc.), the Riemannian Conjugate Gradient (CG) method [2, Ch. 8], was found to be empirically more stable and faster. Below, we provide a short exposition of the CG method in the context of minimizing over \mathbf{B} which belongs to an SPD product manifold.

For an arbitrary non-linear function $\theta(x)$, $x \in \mathbb{R}^n$, the CG method uses the following recurrence at step $k+1$

$$x_{k+1} = x_k + \gamma_k \xi_k, \quad (21)$$

where the direction of descent ξ_k is

$$\xi_k = -\text{grad } \theta(x_k) + \mu_k \xi_{k-1}, \quad (22)$$

with $\text{grad } \theta(x_k)$ defining gradient of θ at x_k ($\xi_0 = -\text{grad } \theta(x_0)$), and μ_k given by

$$\mu_k = \frac{(\text{grad } \theta(x_k))^T (\text{grad } \theta(x_k) - \text{grad } \theta(x_{k-1}))}{\text{grad } \theta(x_{k-1})^T \text{grad } \theta(x_{k-1})}, \quad (23)$$

The step-size γ_k in (21) is usually found via an efficient line-search method [6]. It can be shown that [6, Sec. 1.6] when θ is quadratic with a Hessian \mathbf{Q} , the directions generated by (22) are \mathbf{Q} -conjugate to previous directions of descent $\xi^0, \xi^1, \dots, \xi^{k-1}$; thereby (21) providing the exact minimizer of f in fewer than d iterations (d is the manifold dimension).

For $\mathbf{B} \in \mathbb{M}_n^d$ and referring back to (20), the recurrence in (21) uses the Riemannian retraction [2, Ch. 4] and the gradient $\text{grad}\Theta(\mathbf{B}_k)$ is the Riemannian gradient (here \mathbf{B}_k represents the dictionary tensor at the k -th iteration). This leads to an important issue: the gradients $\text{grad}\Theta(\mathbf{B}_k)$ and $\text{grad}\Theta(\mathbf{B}_{k-1})$ belong to two different tangent spaces $T_{\mathbf{B}_k}\mathbb{M}$ and $T_{\mathbf{B}_{k-1}}\mathbb{M}$ respectively, and thus cannot be combined as in (23). Thus, following [2, Ch. 8] we resort to vector transport – a scheme to transport a tangent vector at $P \in \mathbb{M}$ to a point $\text{Exp}_P(S)$ where $S \in T_P\mathbb{M}$ and Exp is the exponential map. The resulting formula for the direction update becomes

$$\tilde{\xi}_{\mathbf{B}_k} = -\text{grad}\Theta(\mathbf{B}_k) + \mu_k \mathfrak{T}_{\gamma_k \tilde{\xi}_{k-1}}(\tilde{\xi}_{k-1}), \quad (24)$$

where

$$\mu_k = \frac{\langle \text{grad}\Theta(\mathbf{B}_k), \text{grad}\Theta(\mathbf{B}_k) - \mathfrak{T}_{\gamma_k \tilde{\xi}_{k-1}}(\text{grad}\Theta(\mathbf{B}_{k-1})) \rangle}{\langle \text{grad}\Theta(\mathbf{B}_{k-1}), \text{grad}\Theta(\mathbf{B}_{k-1}) \rangle}. \quad (25)$$

Here for $Z_1, Z_2 \in T_P\mathbb{M}$, the map $\mathfrak{T}_{Z_1}(Z_2)$ defines the vector transport given by:

$$\mathfrak{T}_{Z_1}(Z_2) = \left. \frac{d}{dt} \exp_P(Z_1 + tZ_2) \right|_{t=0}. \quad (26)$$

It remains to derive an expression for the Riemannian gradient $\text{grad}\Theta(\mathbf{B})$.

Riemannian Gradient

Lemma 1 connects the Riemannian gradient to the Euclidean gradient of $\Theta(\mathbf{B})$.

Lemma 1. *For a dictionary tensor $\mathbf{B} \in \mathbb{M}_n^d$, let $\Theta(\mathbf{B})$ be a differentiable function. Then the Riemannian gradient $\text{grad}\Theta(\mathbf{B})$ satisfies the following equation:*

$$\langle \text{grad}\Theta(\mathbf{B}), \zeta \rangle_{\mathbf{B}} = \langle \nabla\Theta(\mathbf{B}), \zeta \rangle_I, \forall \zeta \in T_P\mathbb{M}_n^d, \quad (27)$$

where $\nabla\Theta(\mathbf{B})$ is the Euclidean gradient of $\Theta(\mathbf{B})$. The Riemannian gradient for the i -th dictionary atom is given by $\text{grad}_i\Theta(\mathbf{B}) = \mathbf{B}_i \nabla_{\mathbf{B}_i}\Theta(\mathbf{B}) \mathbf{B}_i$.

Proof. See [2, Ch. 5]. The latter expression is obtained by substituting the inner product on the LHS of (27) by its definition in (5).

The Euclidean gradient $\nabla\Theta(\mathbf{B})$ is obtained as follows: let $\mathbf{S}_j = \mathbf{X}_j^{-1/2}$ and $M_j(\mathbf{B}) := \mathbf{B}(\boldsymbol{\alpha}_j) = \sum_{i=1}^n \boldsymbol{\alpha}_j^i \mathbf{B}_i$. Then,

$$\Theta(\mathbf{B}) = \frac{1}{2} \sum_{j=1}^m \text{Tr}(\text{Log}(\mathbf{S}_j M_j(\mathbf{B}) \mathbf{S}_j)^2) + \lambda_{\mathbf{B}} \sum_{i=1}^n \text{Tr}(\mathbf{B}_i). \quad (28)$$

The derivative $\nabla_{\mathbf{B}_i}\Theta(\mathbf{B})$ of (28) w.r.t. to atom \mathbf{B}_i is:

$$\sum_{j=1}^m \boldsymbol{\alpha}_j^i (\mathbf{S}_j \text{Log}(M_j(\mathbf{B})) (M_j(\mathbf{B}))^{-1} \mathbf{S}_j) + \lambda_{\mathbf{B}} I. \quad (29)$$

2.2.2 Sparse Coding Subproblem

Referring back to (19), we now consider the sparse coding subproblem. Given a dictionary tensor \mathbf{B} and a data matrix $\mathbf{X}_j \in \mathbb{S}_+^d$, this subproblem requires solving

$$\begin{aligned} \min_{\boldsymbol{\alpha}_j \geq 0} \quad & \phi(\boldsymbol{\alpha}_j) := \frac{1}{2} d_{\mathbb{R}}^2(\mathbf{X}_j, \mathbf{B}(\boldsymbol{\alpha}_j)) + \text{Sp}(\boldsymbol{\alpha}_j) \\ & = \frac{1}{2} \left\| \text{Log}\left(\sum_{i=1}^n \alpha_j^i \mathbf{X}^{-1/2} \mathbf{B}_j \mathbf{X}^{-1/2}\right) \right\|_F^2 + \text{Sp}(\boldsymbol{\alpha}_j), \end{aligned} \quad (30)$$

where α_j^i is the i -th component of $\boldsymbol{\alpha}_j$ and Sp is a sparsity inducing function. For simplicity, we use $\text{Sp}(\boldsymbol{\alpha}) = \lambda \|\boldsymbol{\alpha}\|_1$, where $\lambda > 0$ is a regularization parameter. Since we are working with $\boldsymbol{\alpha} \geq 0$, we replace this penalty by $\lambda \sum_i \alpha_i$, which is differentiable.

Problem (30) measures reconstruction quality offered by a sparse non-negative linear combination of the atoms to a given input point \mathbf{X} . It will turn out (see experiments in Section 3) that the reconstructions obtained via this model actually lead to significant improvements in performance over sparse coding models that ignore SPD geometry. But this gain comes at a price: objective (30) is difficult to optimize, and remains difficult even if we take into account geodesic convexity of $d_{\mathbb{R}}$.

While in practice this nonconvexity does not seem to hurt our model, we digress below to show a surprising but intuitive constraint under which Problem (30) actually becomes convex. Although we do not exploit this observation to save on computation, we highlight it here due to its theoretical appeal.

Theorem 1 ([14]). *The function $\phi(\boldsymbol{\alpha}) := d_{\mathbb{R}}^2(\sum_i \alpha_i \mathbf{B}_i, \mathbf{X})$ is convex on the set*

$$\mathcal{A} := \{\boldsymbol{\alpha} \mid \sum_i \alpha_i \mathbf{B}_i \preceq \mathbf{X}, \text{ and } \boldsymbol{\alpha} \geq 0\}. \quad (31)$$

2.2.3 Optimizing Sparse Codes

We minimize (30) using a projected gradient method. Specifically, we run the iteration

$$\boldsymbol{\alpha}^{k+1} \leftarrow \mathcal{P}[\boldsymbol{\alpha}^k - \eta_k \nabla \phi(\boldsymbol{\alpha}^k)], \quad k = 0, 1, \dots, \quad (32)$$

where $\mathcal{P}[\cdot]$ denotes the projection operator defined as

$$\mathcal{P}[\boldsymbol{\alpha}] \equiv \boldsymbol{\alpha} \mapsto \operatorname{argmin}_{\boldsymbol{\alpha}'} \frac{1}{2} \|\boldsymbol{\alpha}' - \boldsymbol{\alpha}\|_2^2, \quad \text{s.t.}, \boldsymbol{\alpha}' \in \mathcal{A}. \quad (33)$$

To implement iteration (32) we need to specify three components: (i) the stepsize η_k ; (ii) the gradient $\nabla \phi(\boldsymbol{\alpha}^k)$; and (iii) the projection (33). Lemma 2 shows how to compute the gradient. The projection task (33) is a special least-squares (dual) semidefinite program (SDP), which can be solved using any SDP solver. However, in the interest of speed, we avoid the heavy computational burden imposed by an SDP, and drop the constraint $\boldsymbol{\alpha} \in \mathcal{A}$. Although this sacrifices convexity, the resulting computation is vastly easier, and work well empirically. With this change, we simply have $\mathcal{P}[\boldsymbol{\alpha}] = \max(0, \boldsymbol{\alpha})$.

It remains to specify how to obtain the stepsize η_k . There are several choices in the nonlinear programming literature [6], but most of them can be expensive in our setting. We wish to avoid expensive iterative algorithms for computing η_k , and thus choose to use Barzilai-Borwein stepsizes [5] that have closed forms and that often work remarkably well in practice [5, 43]. In particular, we use the Spectral Projected Gradient (SPG) method [9] by adapting a simplified implementation of [43].

Lemma 2. *Let \mathbf{B} , \mathbf{C} , and \mathbf{X} be fixed SPD matrices. Consider the function $f(x) := \text{d}_R^2(x\mathbf{B} + \mathbf{C}, \mathbf{X})$. The derivative $f'(x)$ is given by*

$$f'(x) = 2 \text{Tr}(\log(\mathbf{S}(x\mathbf{B} + \mathbf{C})\mathbf{S})\mathbf{S}^{-1}(x\mathbf{B} + \mathbf{C})^{-1}\mathbf{B}\mathbf{S}), \quad \text{where } \mathbf{S} = \mathbf{X}^{-1/2}. \quad (34)$$

Proof. Introduce the shorthand $\mathbf{M}(x) \equiv x\mathbf{B} + \mathbf{C}$. Using (5) we have

$$f(x) = \text{Tr}([\log(\mathbf{S}\mathbf{M}(x)\mathbf{S})]^T [\log(\mathbf{S}\mathbf{M}(x)\mathbf{S})]),$$

The chain-rule of calculus then immediately yields the desired result

$$f'(x) = 2 \text{Tr}(\log(\mathbf{S}\mathbf{M}(x)\mathbf{S})(\mathbf{S}\mathbf{M}(x)\mathbf{S})^{-1}\mathbf{S}\mathbf{M}'(x)\mathbf{S}).$$

Writing $\mathbf{M}(\alpha_p) = \alpha_p\mathbf{B}_p + \sum_{i \neq p} \alpha_i\mathbf{B}_i$ and using Lemma 2 we obtain

$$\frac{\partial \phi(\alpha)}{\partial \alpha_p} = \text{Tr}(\log(\mathbf{S}\mathbf{M}(\alpha_p)\mathbf{S})(\mathbf{S}\mathbf{M}(\alpha_p)\mathbf{S})^{-1}\mathbf{S}\mathbf{B}_p\mathbf{S}) + \lambda. \quad (35)$$

Computing (35) for all α is the dominant when running SPG. A naïve implementation of (35) costs $O(nd^3)$, but with slight care this cost can be reduced $O(nd^2) + O(d^3)$ [14].

3 Applications of Sparse Coding

In this section, we describe an application of sparse coding for covariances, namely nearest neighbor (NN) retrieval. This is a fundamental task in several computer vision and machine learning applications in which the goal is to find a data point closest to a given query point within a large database.

3.1 Nearest Neighbors on Covariance Descriptors

Suppose we have obtained a sparse code matrix \mathbf{A} using either GDL or Riemannian DLSC for an input matrix \mathbf{X} . Since, we use an overcomplete dictionary typically only a few dictionary atoms participate in the reconstruction of \mathbf{X} . Thus, with high probability dissimilar input points will obtain different sparse codes. In other words, suppose that we use a dictionary with n rank-one atoms and that only r of these

matrices are used to obtain a reconstruction for a given input. Then, there are $\binom{n}{r}$ unique basis combinations possible. With appropriate choices of n and r , we will likely obtain a unique set of rank-one matrices that encode \mathbf{X} .

Using this intuition, we propose a sorted integer tuple representation to encode an input covariance matrix; the integers simply index the dictionary atoms used in the encoding. Formally, let $\mathbf{X} \in \mathbb{S}_+^d$ be the input, \mathbf{B} an overcomplete dictionary, and \mathbf{u}_i ($i \in [n]$) a unique identifier for the i -th atom of \mathbf{B} . If $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$ is the coefficient vector corresponding to $\mathbf{X} \approx \mathbf{B}(\boldsymbol{\alpha})$, then the tuple $h(\mathbf{X}) = \langle \mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_k} \rangle$ is the *hash code* of \mathbf{X} . We choose only those identifiers for which the code α_j is larger than a threshold $\varepsilon \geq 0$.

In our case, we assume that the u_i 's are just integers in $\{1, \dots, n\}$ and that the hash code is a sorted tuple of these indices. The threshold ε helps select significant coefficients from the sparse coding, and makes the chosen code robust to noise. This coded representation enables the use of hash tables for fast locality sensitive hashing. Let us see how. Each column of the dictionary is identified by its index number; so each hash-key is a set of integers encoded as a character string. To tackle collisions in the hash buckets, the colliding input matrices are organized as a linked list. If the linked list gets too long, the data within a hash bucket can be further organized using a metric tree or any other efficient data structure. This idea of hashing is a direct adaptation of the scheme proposed in [12].

Given a query SPD matrix, we solve the sparse coding problem to first obtain the corresponding sparse coefficients. From these coefficients we compute the above hash code query the hash table. If there are several entries in a matching bucket, we run a linear scan using the AIRM distance (5) to find the best matches (the bucket can also be organized for faster than linear scans, if desired).

3.2 GDL Experiments

This section illustrates nearest neighbor search based upon our dictionary learning examples. We use the following datasets:

- **Face recognition.** The *FERET face dataset* [42, 41] contains facial appearances segregated into multiple classes. Each class has different views of the face of the same person for varying poses. We selected six images from each class. Inspired by the success of covariances created from Gabor filters for face recognition [31], we applied 40 Gabor filters on each image, later combining the filters into a covariance of size 40×40 . We created a covariance dataset of approximately 10K descriptors using this approach.
- **Texture classification.** Texture is an essential cue in many data mining applications like satellite imagery, industry inspection systems, etc. Thus, we used a combination of the *Brodatz dataset* [11] and the *Curret dataset* [16] for creating a texture covariance dataset. Brodatz dataset contains approximately 111 texture classes, while Curret dataset contains 60 classes. To create the covariances data, we used the feature vector $F = [x, y, I, I_x, I_y]$, where the first two dimensions are

the relative location of the pixel with respect to the texture patch, the third dimension encodes the grayscale intensity, and the last two dimensions capture the pixel gradients. Thus, each covariance is 5×5 , and we created approximately 40K such covariances.

Methods Compared. We compare against locality sensitive hashing (LSH) of vectorized covariances (VEC), hashing after log-Euclidean embedding (L2LSH), and kernelized LSH [28] using an RBF kernel using the AIRM distance.

3.2.1 GDL Experimental Setup

We implemented GDL in Matlab. For L2LSH, VEC, and HAM we used the C-implementation from the Caltech Toolbox.¹ Since the programs have different computational baselines, we cannot compare their retrieval speed. Rather, we show in Table 1 the average portion of each of the datasets scanned by GDL to find the nearest neighbor. The geodesic distance was used to resolve hash table collisions. As is seen from the table, the percentage coverage is low, which is exactly as desired.

Dataset	Faces	Texture
Avg. coverage (%)	3.54	6.26

Table 1 Percentage of the database searched to find the nearest neighbor using sparse codes generated by GDL.

Next, we substantiate effectiveness of our NN retrieval scheme. To this end, we split each of the datasets into database and query sets (approximately 5% of the data). To compute the ground truth we use a linear scan over the database using geodesic distances to measure nearness. Since ensuring exact NN is hard, we restrict our search to Approximate Nearest Neighbors (ANN). Assume \mathbf{Q} is a query point, \mathbf{X}_{ls} is the exact NN found by a linear scan and \mathbf{X}_{algo} is the neighbor returned by an NN algorithm. We classify an NN as correct if $\frac{d_{AIRM}(\mathbf{Q}, \mathbf{X}_{ls})}{d_{AIRM}(\mathbf{Q}, \mathbf{X}_{AIRM})} > \epsilon$; we use $\epsilon = 0.75$ below. Fig. 2 shows the accuracy of different methods, where

$$\text{Accuracy} := \frac{\#\text{correct matches}}{\#\text{query size}}. \quad (36)$$

The plots indicate that GDL performs well across the datasets, while performance of the other methods varies. Vectorizing input matrices fails on all datasets, while KLSH performs reasonably well. We note, however, that KLSH needs to compute the kernel matrix for the query point against the *entire* dataset—this can drastically slow it down. On the face dataset, all methods had high accuracy, most probably because this dataset is noise free.

¹ <http://www.vision.caltech.edu/malaa/software/research/image-search/>

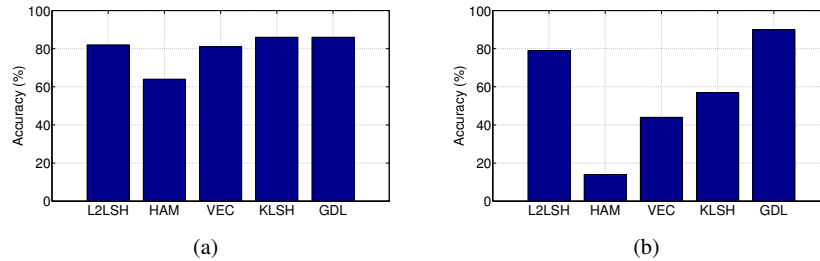


Fig. 2 Plots demonstrating nearest neighbor classification accuracy of GDL generated sparse codes compared to various standard techniques; (a) on faces dataset and (b) on texture dataset.

3.3 Riemannian Dictionary Learning Experiments

Next, we evaluate the Riemannian DLSC setup, and denote below dictionary learning by DL and sparse coding by SC. We compare our Riemannian (Riem) formulation against combinations of several state-of-the-art DLSC methods on SPD matrices, namely (i) log-Euclidean (LE) metric for DLSC [25], (ii) Frobenius norm (Frob) which discards the manifold structure, (iii) kernel methods such as the Stein-Kernel [48] proposed in [26], and the log-Euclidean kernel [30].

We experiment on data available from three standard computer vision applications: (i) 3D object recognition on the RGBD objects dataset [29]; (ii) texture recognition on the standard Brodatz dataset [38]; and (iii) person re-identification on the ETHZ people dataset [19]. We describe (i) and (iii) below.

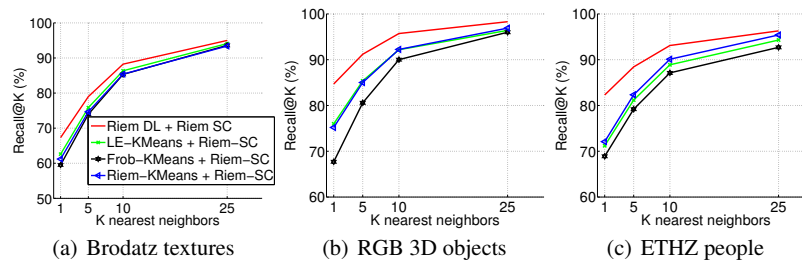


Fig. 3 Results of Nearest neighbor recall@K accuracy against increasing number of retrieved points (K). Comparisons of Riem-DL and Riem-SC against other DL learning schemes based on clustering, while using our Riem-SC for sparse coding.

- **Person re-identification task.** We use the benchmark ETHZ dataset [44] for evaluating people re-identification. This dataset consists of low-resolution images of tracked people from a real-world surveillance setup. The images are from 146 different individuals. There are about 5–356 images per person. There are a total

of 8580 images in this dataset. Rather than detailing the results on several feature combinations, we describe here the feature combination that worked the best in our experiments. For this purpose, we used a validation set of 500 covariances and 10 true clusters from this dataset. The performance was evaluated using the Log-Euclidean SC setup with a dictionary learning via Log-Euclidean K-Means. We used a combination of nine features for each image as described below:

$$F_{ETHZ} = [x, I_r, I_g, I_b, Y_i, |I_x|, |I_y|, |\sin(\theta) + \cos(\theta)|, |H_y|],$$

where x is the x-coordinate of a pixel location, I_r, I_g, I_b are the RGB color of a pixel, Y_i is the pixel intensity in the YCbCr color space, I_x, I_y are the gray scale pixel gradients, and H_y is the y-gradient of pixel hue. We also use the gradient angle $\theta = \tan^{-1}(I_y/I_x)$ in our feature set. Each image is resized to a fixed size 300×100 , and divided into upper and lower parts. We compute two different region covariances for each part, which are combined as two block diagonal matrices to form a single covariance of size 18×18 for each appearance image.

- **3D Object Recognition.** The goal of this experiment is to recognize objects in 3D point clouds. We use the public RGB-D Object dataset [29], which consists of about 300 objects belonging to 51 categories and spread across $\sim 250K$ frames. We used approximately 15K frames for our evaluation with approximately 250–350 frames devoted to every object seen from three different view-points (30, 45, and 60 degrees above the horizon). Following the procedure suggested in [23][Chap. 5], for every frame, the object was segmented out and 18 dimensional feature vectors generated for every 3D point in the cloud (and thus 18×18 covariance descriptors); the features we used are as follows:

$$F_{RGBD} = [x, y, z, I_r, I_g, I_b, I_x, I_y, I_{xx}, I_{yy}, I_{xy}, I_m, \delta_x, \delta_y, \delta_m, v_x, v_y, v_z], \quad (37)$$

where the first three dimensions are the spatial coordinates, I_m is the magnitude of the intensity gradient, δ 's represent gradients over the depth-maps, and v represents the surface normal at the given 3D point.

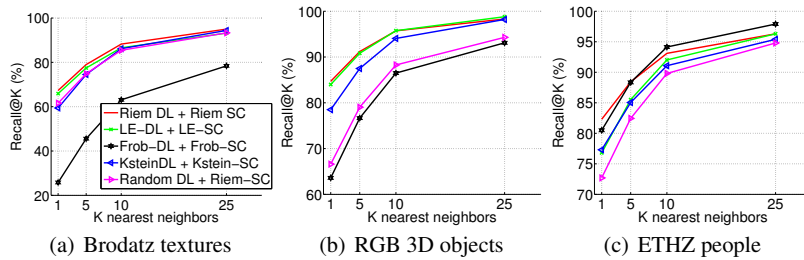


Fig. 4 Results of k-nearest neighbor retrieval accuracy against an increasing number of retrieved points (K). Comparisons of Riem-DL and Riem-SC against other dictionary learning and sparse coding combinations.

3.3.1 Evaluation Techniques

We evaluate our algorithms for nearest neighbor (NN) retrieval against a gallery set via computing the Euclidean distances between sparse codes. We use the standard Recall@K accuracy defined as follows: Given a gallery \mathcal{X} and a query set \mathcal{Q} . Recall@K computes the average accuracy when retrieving K nearest neighbors from \mathcal{X} for each instance in \mathcal{Q} . Suppose G_K^q is the set of ground truth class labels associated with the q th query, and S_K^q is the set of labels associated with the K neighbors found by some algorithm, then

$$\text{Recall@K} := \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \frac{|G_K^q \cap S_K^q|}{|G_K^q|}. \quad (38)$$

All the experiments used 5-fold cross-validation in which 80% of the datasets were used for training the dictionary, 10% for generating the gallery set, and the rest as queries. The size of the dictionary was considered to be twice the number of classes in the respective dataset. This scheme was considered for all the comparison methods as well. Our DLSC scheme was implemented in MATLAB. We used the Manopt optimization toolbox [10] for implementing the CG method for our DL subproblem. We found that initializing the dictionary learning setup using K-Means clustering (using the Karcher mean algorithm [40]) led to faster convergence of CG.

3.3.2 Results

We compare below performance of our Riem-DL and Riem-SC against several prior DLSC schemes on the three datasets described above. In particular, we compare (i) Riemannian geometric methods such as log-Euclidean (LE-DL + LE-SC), (ii) Kernelized methods using the Stein kernel (Kernel-Stein-DL and kernel-Stein-SC), (iii) Euclidean DLSC (Frob-DL + Frob-SC), and using a dictionary generated by random sampling the dataset followed by sparse coding using our Riemannian method (Random-DL + Riem-SC). In Figure 4, we show the performance comparison for the task of K-NN where K is increased from 1 to 25.

A commonly adopted alternative to dictionary learning is to approximate the dictionary using centroids of a K-Means clustering of the dataset. Such a method is faster than Riemannian DL, and also demonstrate reasonable performance [13, 46]. Thus, an important experiment is to ensure that learning the dictionary actually provides superior performance compared to the *ad hoc* clustering setup. In Figure 3, we plot the K-NN retrieval when we use a clustering scheme to generate the dictionary.

3.3.3 Discussion of Results

With regard to Figure 4, we found that the performance of different methods is diverse across datasets. For example, the log-euclidean DLSC variant (LE-DL+LE-

SC) is generally seen to show good performance, but its performance is inferior when the number of data instances per class is small (as in the ETHZ people dataset). The kernelized DLSC method (Kernel-Stein-DL) performs favorably on most datasets. The most surprising of the results that we found was for Frob-DL. It is generally assumed that using Frobenius distance for comparing SPD matrices leads to poor accuracy, a view echoed by Figures 4(a) and 4(b). However, when the matrices are ill-conditioned, taking the logarithm (as in the LE-DL scheme) of these matrices results in amplifying the influence of the smaller eigenvalues, which is essentially noise. When learning a dictionary, the atoms will be learned to reconstruct this noise against the signal, thus leading to inferior performance than for FrobDL which do not use the logarithm. In comparison to all the compared methods, Riem-DL+Riem-SC was found to produce consistent and competitive performance, substantiating the usefulness of our model. While running the experiments, we found that the initialization of our DL sub-problem (using Riemannian K-Means) played an important role in achieving this superior performance.

We further compare Riem-DL against alternative DL schemes via clustering in Figure 3. We see that learning the dictionary using Riem-DL demonstrates the best performance against the next best and efficient alternative of using LE-KMeans as was done in [13]. Using Frob-KMeans or using a random dictionary are generally seen to have inferior performance compared to other learning methods.

3.4 GDL versus Riemannian Sparse Coding

Finally, we compare sparse coding via our GDL model and the Riem-SC setup. We use the Brodatz and RGB-D Object recognition datasets. Tables 2 and 3 show the results. While GDL is significantly faster in sparse coding, as is clear from the table, the Riemannian approach leads to much higher accuracy.

Method	Accuracy (%)
Frob-SC	32.3 (4.4)
TSC [46]	35.6 (7.1)
GDL	43.7 (6.3)
Riem-SC	53.9 (3.4)

Table 2 Brodatz texture dataset

Comparison of the average classification accuracy using a linear SVM.

Method	Accuracy (%)
Frob-SC	80.3 (1.1)
TSC [46]	72.8 (2.1)
GDL	61.9 (0.4)
Riem-SC	84.0 (0.6)

Table 3 RGB-D Object Recognition

4 Conclusion and Future Work

In this chapter, we reviewed the steps for constructing covariance descriptors, followed by a brief exposition of SPD matrix geometry motivated by the design of

novel machine learning models on these descriptors. We covered the concrete problems of dictionary learning and sparse coding, and noted two approaches: (i) a framework that uses Euclidean embedding of SPD matrices for sparse coding; and (ii) a Riemannian geometric approach. Our experiments demonstrated that designing machine learning algorithms for SPD matrices that respect the Riemannian geometry fares significantly better than using Euclidean embedding.

That said, Riemannian optimization algorithms are usually computationally more expensive. This is mainly due to the need for operating in the tangent space of the SPD manifold, which involves matrix exponentials and logarithms that require $O(d^3)$ flops. Designing faster Riemannian machine learning algorithms is a challenge that needs to be addressed for these algorithms to be more widely accepted.

Acknowledgements AC is funded by the Australian Research Council Centre of Excellence for Robotic Vision (number CE140100016). SS acknowledges support from NSF grant IIS-1409802.

References

1. Absil, P.A., Baker, C.G., Gallivan, K.A.: Trust-region methods on riemannian manifolds. *Foundations of Computational Mathematics* **7**(3), 303–330 (2007)
2. Absil, P.A., Mahony, R., Sepulchre, R.: *Optimization algorithms on matrix manifolds*. Princeton University Press (2009)
3. Alexander, D.C., Pierpaoli, C., Basser, P.J., Gee, J.C.: Spatial transformations of diffusion tensor magnetic resonance images. *IEEE Transactions on Medical Imaging* **20**(11), 1131–1139 (2001)
4. Arsigny, V., Fillard, P., Pennec, X., Ayache, N.: Log-Euclidean metrics for fast and simple calculus on diffusion tensors. *Magnetic Resonance in Medicine* **56**(2), 411–421 (2006)
5. Barzilai, J., Borwein, J.M.: Two-Point Step Size Gradient Methods. *IMA J. Num. Anal.* **8**(1) (1988)
6. Bertsekas, D.P., Bertsekas, D.P.: *Nonlinear Programming*, second edn. Athena Scientific (1999)
7. Bhatia, R.: *Positive Definite Matrices*. Princeton University Press (2007)
8. Birgin, E., Martínez, J., Raydan, M.: Nonmonotone spectral projected gradient methods on convex sets. *SIAM Journal on Optimization* **10**(4), 1196–1211 (2000)
9. Birgin, E.G., Martínez, J.M., Raydan, M.: Algorithm 813: SPG - Software for Convex-constrained Optimization. *ACM Transactions on Mathematical Software* **27**, 340–349 (2001)
10. Boumal, N., Mishra, B., Absil, P.A., Sepulchre, R.: Manopt, a matlab toolbox for optimization on manifolds. *The Journal of Machine Learning Research* **15**(1), 1455–1459 (2014)
11. Brodatz, P.: *Textures: a photographic album for artists and designers*, vol. 66. Dover New York (1966)
12. Cherian, A.: Nearest neighbors using compact sparse codes. In: *International Conference on Machine Learning*, pp. 1053–1061 (2014)
13. Cherian, A., Sra, S.: Riemannian sparse coding for positive definite matrices. In: *European Conference on Computer Vision*. Springer (2014)
14. Cherian, A., Sra, S.: Riemannian dictionary learning and sparse coding for positive definite matrices. *arXiv preprint arXiv:1507.02772* (2015)
15. Cherian, A., Sra, S., Banerjee, A., Papanikolopoulos, N.: Jensen-bregman logdet divergence with application to efficient similarity search for covariance matrices. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **35**(9), 2161–2174 (2013)

16. Dana, K., Van Ginneken, B., Nayar, S., Koenderink, J.: Reflectance and texture of real-world surfaces. *ACM Transactions on Graphics (TOG)* **18**(1), 1–34 (1999)
17. Doderer, L., Minh, H.Q., Biagio, M.S., Murino, V., Sona, D.: Kernel-based classification for brain connectivity graphs on the Riemannian manifold of positive definite matrices. In: *International Symposium on Biomedical Imaging*, pp. 42–45. IEEE (2015)
18. Elad, M., Aharon, M.: Image denoising via learned dictionaries and sparse representation. In: *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, vol. 1, pp. 895–900. IEEE (2006)
19. Ess, A., Leibe, B., Gool, L.V.: Depth and appearance for mobile scene analysis. In: *International Conference on Computer Vision*. IEEE (2007)
20. Fano, U.: Description of states in quantum mechanics by density matrix and operator techniques. *Reviews of Modern Physics* **29**(1), 74–93 (1957)
21. Faraki, M., Harandi, M.: Bag of riemannian words for virus classification. *Case Studies in Intelligent Computing: Achievements and Trends* pp. 271–284 (2014)
22. Fehr, D., Cherian, A., Sivalingam, R., Nickolay, S., Morellas, V., Papanikolopoulos, N.: Compact covariance descriptors in 3d point clouds for object recognition. In: *IEEE International Conference on Robotics and Automation*, pp. 1793–1798. IEEE (2012)
23. Fehr, D.A.: Covariance based point cloud descriptors for object detection and classification. University Of Minnesota (2013)
24. Ferro-Famil, L., Pottier, E., Lee, J.: Unsupervised classification of multifrequency and fully polarimetric SAR images based on the H/A/Alpha-Wishart classifier. *IEEE Transactions on Geoscience and Remote Sensing* **39**(11), 2332–2342 (2001)
25. Guo, K., Ishwar, P., Konrad, J.: Action recognition using sparse representation on covariance manifolds of optical flow. In: *Advanced Video and Signal Based Surveillance*. IEEE (2010)
26. Harandi, M.T., Sanderson, C., Hartley, R., Lovell, B.C.: Sparse coding and dictionary learning for symmetric positive definite matrices: A kernel approach. In: *European Conference on Computer Vision*. Springer (2012)
27. Jayasumana, S., Hartley, R., Salzmann, M., Li, H., Harandi, M.: Kernel methods on the Riemannian manifold of symmetric positive definite matrices. In: *Computer Vision and Pattern Recognition*. IEEE (2013)
28. Kulis, B., Grauman, K.: Kernelized locality-sensitive hashing for scalable image search. *ICCV*, October **1**, 3 (2009)
29. Lai, K., Bo, L., Ren, X., Fox, D.: A large-scale hierarchical multi-view RGB-D object dataset. In: *International Conference on Robotics and Automation* (2011)
30. Li, P., Wang, Q., Zuo, W., Zhang, L.: Log-euclidean kernels for sparse representation and dictionary learning. In: *International Conference on Computer Vision*. IEEE (2013)
31. Liu, C.: Gabor-based kernel PCA with fractional power polynomial models for face recognition. *Pattern Analysis and Machine Intelligence* **26**(5), 572–581 (2004)
32. Liu, D., Nocedal, J.: On the limited memory BFGS method for large scale optimization. *Mathematical programming* **45**(1), 503–528 (1989)
33. Ma, B., Su, Y., Jurie, F.: BiCov: a novel image representation for person re-identification and face verification. In: *British Machine Vision Conference* (2012)
34. Mairal, J., Bach, F., Ponce, J.: Sparse modeling for image and vision processing. *arXiv preprint arXiv:1411.3230* (2014)
35. Marčelja, S.: Mathematical description of the responses of simple cortical cells*. *JOSA* **70**(11), 1297–1300 (1980)
36. O. Tuzel, F. Porikli, and P. Meer.: Covariance Tracking using Model Update Based on Lie Algebra. *Computer Vision and Pattern Recognition* (2006)
37. O. Tuzel, F. Porikli, and P. Meer.: Region Covariance: A Fast Descriptor for Detection and Classification. In: *European Conference on Computer Vision* (2006)
38. Ojala, T., Pietikäinen, M., Harwood, D.: A comparative study of texture measures with classification based on featured distributions. *Pattern recognition* **29**(1), 51–59 (1996)
39. Pang, Y., Yuan, Y., Li, X.: Gabor-based region covariance matrices for face recognition. *IEEE Transactions on Circuits and Systems for Video Technology* **18**(7), 989–993 (2008)

40. Pennec, X., Fillard, P., Ayache, N.: A Riemannian framework for tensor computing. *International Journal of Computer Vision* **66**(1), 41–66 (2006)
41. Phillips, P., Moon, H., Rizvi, S., Rauss, P.: The FERET evaluation methodology for face-recognition algorithms. *Pattern Analysis and Machine Intelligence* **22**(10), 1090–1104 (2000)
42. Phillips, P., Wechsler, H., Huang, J., Rauss, P.: The FERET database and evaluation procedure for face-recognition algorithms. *Image and Vision Computing* **16**(5), 295–306 (1998)
43. Schmidt, M., van den Berg, E., Friedlander, M., Murphy, K.: Optimizing Costly Functions with Simple Constraints: A Limited-Memory Projected Quasi-Newton Algorithm. In: *International Conference on Artificial Intelligence and Statistics* (2009)
44. Schwartz, W., Davis, L.: Learning Discriminative Appearance-Based Models Using Partial Least Squares. In: *Proceedings of the XXII Brazilian Symposium on Computer Graphics and Image Processing* (2009)
45. Shinohara, Y., Masuko, T., Akamine, M.: Covariance clustering on Riemannian manifolds for acoustic model compression. In: *International Conference on Acoustics, Speech and Signal Processing* (2010)
46. Sivalingam, R., Boley, D., Morellas, V., Papanikolopoulos, N.: Tensor sparse coding for region covariances. In: *European Conference on Computer Vision*. Springer (2010)
47. Somasundaram, G., Cherian, A., Morellas, V., Papanikolopoulos, N.: Action recognition using global spatio-temporal features derived from sparse representations. *Computer Vision and Image Understanding* **123**, 1–13 (2014)
48. Sra, S.: Positive Definite Matrices and the S-Divergence. *Proceedings of the American Mathematical Society* (2015). ArXiv:1110.1773v4
49. Sra, S., Cherian, A.: Generalized dictionary learning for symmetric positive definite matrices with application to nearest neighbor retrieval. In: *European Conference on Machine Learning*. Springer (2011)
50. Su, J., Srivastava, A., de Souza, F., Sarkar, S.: Rate-invariant analysis of trajectories on riemannian manifolds with application in visual speech recognition. In: *Computer Vision and Pattern Recognition*, pp. 620–627. IEEE (2014)
51. Tosato, D., Farenzena, M., Spera, M., Murino, V., Cristani, M.: Multi-class classification on Riemannian manifolds for video surveillance. In: *European Conference on Computer Vision* (2010)
52. Wang, R., Guo, H., Davis, L.S., Dai, Q.: Covariance discriminative learning: A natural and efficient approach to image set classification. In: *Computer Vision and Pattern Recognition*. IEEE (2012)
53. Zheng, W., Tang, H., Lin, Z., Huang, T.S.: Emotion recognition from arbitrary view facial images. In: *European Conference on Computer Vision*, pp. 490–503. Springer (2010)