

Learning to Optimize on SPD Manifolds

Supplementary Materials

Zhi Gao¹, Yuwei Wu^{1*}, Yunde Jia¹, Mehrtash Harandi²

¹Beijing Laboratory of Intelligent Information Technology

School of Computer Science, Beijing Institute of Technology, Beijing, China

²Department of Electrical and Computer Systems Eng., Monash University, and Data61, Australia

{gaozhi_2017, wuyuwei, jiayunde}@bit.edu.cn, mehrtash.harandi@monash.edu

1. Training Our Optimizer on Specific Tasks

Recall that optimization problems with SPD constraints can be formulated as minimizing a loss

$$\mathcal{L}(\mathbf{M}) \triangleq \frac{1}{n} \sum_{i=1}^n l(f(\mathbf{M}, \mathbf{x}_i), \mathbf{y}_i), \quad (1)$$

with respect to \mathbf{M} , that is $\min_{\mathbf{M} \in \mathcal{S}_{++}^d} \mathcal{L}(\mathbf{M})$. In Eq. (1), \mathbf{x}_i is the i -th training sample, and its corresponding target (e.g., label) is \mathbf{y}_i . $f(\cdot)$ and $l(\cdot)$ represent the prediction and objective functions, respectively, and $\mathbf{M} \in \mathcal{S}_{++}^d$ encapsulates parameters of the optimization problem.

To solve Eq. (1), we parameter the SPD optimizer by a network, and ϕ is the parameter of the network, through which the SPD parameter is updated by

$$\mathbf{M}^{(t+1)} = \Gamma_{\mathbf{M}^{(t)}}(-g_\phi(\nabla_{\mathbf{M}}^{(t)} \mathbf{S}^{(t-1)})), \quad (2)$$

where $\Gamma_{\mathbf{M}^{(t)}}(\cdot)$ is the retraction operation, and $g_\phi(\nabla_{\mathbf{M}}^{(t)} \mathbf{S}^{(t-1)})$ is the update vector on the tangent space. We train the optimizer by minimizing the following meta-objective

$$\begin{aligned} \mathcal{J}(\phi) &= \frac{1}{m} \sum_t^T \sum_j^m \mathcal{L}(\mathbf{M}_j^{(t+1)}) \\ &= \frac{1}{mn} \sum_t^T \sum_{j,i}^{m,n} l\left(f\left(\mathbf{M}_j^{(t+1)}, \mathbf{x}_i\right), \mathbf{y}_i\right) \\ &= \frac{1}{mn} \sum_t^T \sum_{j,i}^{m,n} l\left(f\left(\Gamma_{\mathbf{M}_j^{(t)}}\left(-g_\phi\left(\nabla_{\mathbf{M}_j}^{(t)} \mathbf{S}_j^{(t-1)}\right)\right), \mathbf{x}_i\right), \mathbf{y}_i\right), \end{aligned} \quad (3)$$

where m is the batchsize in the outer loop, n is the batchsize in the inner loop, and we consider T continuous steps in the inner loop once. Note that, forms of the projection function $f(\cdot)$ and the objective function $l(\cdot)$ differ per task. In this section, we will elaborate on specific forms of $f(\cdot)$ and $l(\cdot)$ on the metric nearness, clusters, and similarity learning tasks.

*Corresponding author

1.1. Metric Nearness

Metric nearness refers to the problem of optimally restoring metric properties to distance measurements that happen to be non-metric due to measurement errors or otherwise [1]. In the metric nearness task, given a set of vectors $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$ and an SPD matrix $\mathbf{A} \in \mathcal{S}_{++}^d$, we expect that we can discover another SPD parameter \mathbf{M} to project the feature $\mathbf{A}\mathbf{x}_i$ back to its original value \mathbf{x}_i in the metric space. The prediction function is $f(\mathbf{M}, \mathbf{x}_i) = \mathbf{M}\mathbf{A}\mathbf{x}_i$, and the objective function is $l(f(\mathbf{M}, \mathbf{x}_i), \mathbf{x}_i) = \|\mathbf{M}\mathbf{A}\mathbf{x}_i - \mathbf{x}_i\|_2^2$, where $\|\cdot\|_2$ is the 2-norm. Thus, the metric nearness task is formulated as minimizing the loss

$$\min_{\mathbf{M} \in \mathcal{S}_{++}^d} \mathcal{L}(\mathbf{M}) \triangleq \frac{1}{n} \sum_i^n \|\mathbf{M}\mathbf{A}\mathbf{x}_i - \mathbf{x}_i\|_2^2. \quad (4)$$

To train our optimizer, we build the meta-objective as

$$\begin{aligned} \min_{\phi} \mathcal{J}(\phi) &= \frac{1}{m} \sum_t^T \sum_j^m \mathcal{L}(\mathbf{M}_j^{(t+1)}) \\ &= \frac{1}{mn} \sum_{t,j,i}^{T,m,n} \|\mathbf{M}_j^{(t+1)} \mathbf{A}\mathbf{x}_i - \mathbf{x}_i\|_2^2 \\ &= \frac{1}{mn} \sum_{t,j,i}^{T,m,n} \|\Gamma_{\mathbf{M}_j^{(t)}}(-g_\phi(\nabla_{\mathbf{M}_j}^{(t)} \mathbf{S}_j^{(t-1)})) \mathbf{A}\mathbf{x}_i - \mathbf{x}_i\|_2^2. \end{aligned} \quad (5)$$

Updating ϕ . In the outer loop, we resort to the gradient descent to minimize Eq. (5), where the derivative $\frac{\partial \mathcal{J}(\phi)}{\partial \phi}$ is needed. We calculate $\frac{\partial \mathcal{J}(\phi)}{\partial \phi}$ by the truncated backpropagation through time (TBPTT) algorithm. As the meta-objective contains T steps in the inner loop, we denote the loss at step t as $\mathcal{J}^{(t)}(\phi) = \frac{1}{m} \sum_j^m \mathcal{L}(\mathbf{M}_j^{(t+1)})$, and $\mathcal{J}(\phi) = \sum_t^T \mathcal{J}^{(t)}(\phi)$. We first compute $\frac{\partial \mathcal{J}^{(t)}(\phi)}{\partial \phi}$ for each step t . Then we compute $\frac{\partial \mathcal{J}(\phi)}{\partial \phi}$ by $\frac{\partial \mathcal{J}(\phi)}{\partial \phi} = \sum_t^T \frac{\partial \mathcal{J}^{(t)}(\phi)}{\partial \phi}$. In this process, the derivatives of the update vector $\mathbf{P}^{(t)}$ and

the SPD parameter $M^{(t)}$ with respect to $\mathcal{J}^{(t)}(\phi)$ are non-trivial in the retraction operation $\Gamma_{M_j^{(t)}}(\cdot)$, as $\Gamma_{M_j^{(t)}}(\cdot)$ contains the matrix power and matrix exponential operations: $M^{(t)\frac{1}{2}}$, $M^{(t)\frac{-1}{2}}$, and $\expm(-M^{(t)\frac{-1}{2}}P^{(t)}M^{(t)\frac{1}{2}})$. To solve this issue, we rewrite the retraction operation as

$$\begin{cases} U\Sigma U^\top = M^{(t)} \\ U_Q \Sigma_Q U_Q^\top = Q = -U\Sigma^{\frac{-1}{2}}U^\top P^{(t)}U\Sigma^{\frac{-1}{2}}U^\top \\ M^{(t+1)} = (U\Sigma^{\frac{1}{2}}U^\top)(U_Q \exp(\Sigma_Q)U_Q^\top)(U\Sigma^{\frac{1}{2}}U^\top) \end{cases}, \quad (6)$$

where $U\Sigma U^\top = M^{(t)}$ and $U_Q \Sigma_Q U_Q^\top = Q$ are the eigenvalue decomposition. After obtaining $\frac{\partial \mathcal{J}}{\partial U_Q}$ and $\frac{\partial \mathcal{J}}{\partial \Sigma_Q}$, we can compute $\frac{\partial \mathcal{J}}{\partial Q}$ based on Proposition 1, and $\frac{\partial \mathcal{J}}{\partial P^{(t)}}$ can be computed by $\frac{\partial \mathcal{J}}{\partial P^{(t)}} = -U\Sigma^{\frac{-1}{2}}U^\top \frac{\partial \mathcal{J}}{\partial Q} U\Sigma^{\frac{-1}{2}}U^\top$. Similarly, after obtaining $\frac{\partial \mathcal{J}}{\partial U}$ and $\frac{\partial \mathcal{J}}{\partial \Sigma}$, we can compute $\frac{\partial \mathcal{J}}{\partial M^{(t)}}$. Thus, $\frac{\partial \mathcal{J}^{(t)}(\phi)}{\partial \phi}$ can be computed by the channel rule according to our optimizer architecture.

Algorithm 1 Training Process of Our Optimizer

Input: The randomly initialized optimizer parameter. The initial SPD parameter $M^{(0)} = \mathbf{I}_d$. The initial state $\mathbf{S}^{(0)} = \mathbf{0}_d$. The initial experience pool $\Psi = \emptyset$.

Output: The optimizer parameter ϕ .

while not reach the maximum iteration of the observation stage **do**

Compute the loss \mathcal{L} of the base learner by Eq. (4);
 Compute the gradient $\nabla_{M_j}^{(t)}$;
 Update the parameter $M^{(t+1)}$ by Eq. (2);
 Push $\{M^{(t+1)}, \mathbf{S}^{(t)}\}$ into Ψ ;

end

while not reach the maximum iteration of the learning stage **do**

Randomly select $\{(M_j^{(t)}, \mathbf{S}_j^{(t-1)})\}_{j=1}^m$ from Ψ ;
while not reach T **do**
 Compute the loss \mathcal{L} of the base learner by Eq. (4);
 Compute the gradient $\nabla_{M_j}^{(t)}$;
 Update $M^{(t+1)}$ by Eq. (2);
end

Compute the loss \mathcal{J} of our optimizer by Eq. (5);

while not reach T **do**

Compute $\frac{\partial \mathcal{J}^{(t)}(\phi)}{\partial \phi}$;

end

Compute $\frac{\partial \mathcal{J}(\phi)}{\partial \phi} = \sum_t^T \frac{\partial \mathcal{J}^{(t)}(\phi)}{\partial \phi}$;

Update ϕ via the ADAM algorithm;

if $t + T > \tau$ **then**

Push $\{(M_j^{(0)} = \mathbf{I}_d, \mathbf{S}_j^{(0)} = \mathbf{0}_d)\}_{j=1}^m$ into Ψ ;

else

Push $\{(M_j^{(t+T)}, \mathbf{S}_j^{(t+T-1)})\}_{j=1}^m$ into Ψ ;

end

end

Return the parameter ϕ of our optimizer.

Since the derivative $\frac{\partial \mathcal{J}(\phi)}{\partial \phi} = \sum_t^T \frac{\partial \mathcal{J}^{(t)}(\phi)}{\partial \phi}$ is obtained, we utilized the ADAM algorithm to update ϕ , and the learn-

ing rate was set as 0.001. A more detailed training process of our optimizer is shown in Algorithm 1. After obtaining the parameter ϕ of our optimizer, our optimizer can work well on the metric nearness task.

Proposition 1. Let $U\Sigma U^\top = M$ be the eigenvalue decomposition of M . For the loss function $\mathcal{J}(U, \Sigma)$, given the derivatives $\frac{\partial \mathcal{J}}{\partial U}$ and $\frac{\partial \mathcal{J}}{\partial \Sigma}$, the derivative $\frac{\partial \mathcal{J}}{\partial M}$ is

$$\frac{\partial \mathcal{J}}{\partial M} = 2U(R^\top \otimes (U^\top \frac{\partial \mathcal{J}}{\partial U} + \frac{\partial \mathcal{J}}{\partial U}^\top U))U^\top + U(\frac{\partial \mathcal{J}}{\partial \Sigma})_{diag}U^\top, \quad (7)$$

where

$$R_{ij} = \begin{cases} \frac{1}{\lambda_i - \lambda_j}, & \text{if } i \neq j \\ 0, & \text{if } i = j \end{cases},$$

λ_i is the i -th eigenvalue, and \mathbf{X}_{diag} denotes \mathbf{X} with all off-diagonal elements being 0.

1.2. Clustering

In our visual clustering task, each image is represented by a covariance descriptor [2, 3, 5], and we expect to divide all covariance descriptors into a number of clusters and find their centers. For each image, we resize it to 128×128 and divide the image into 1024 grids. The size of each grid is 4×4 . At the position (u, v) of the grid, we extract a 5-D feature vector as

$$\mathbf{x}_{uv} = \left[I_{uv}, \left| \frac{\partial I}{\partial u} \right|, \left| \frac{\partial I}{\partial v} \right|, \left| \frac{\partial I^2}{\partial u^2} \right|, \left| \frac{\partial I^2}{\partial v^2} \right| \right]. \quad (8)$$

Based on extracted features, we compute a 5×5 covariance descriptor to represent the image:

$$\mathbf{X} = \sum_{uv} (\mathbf{x}_{uv} - \boldsymbol{\mu})(\mathbf{x}_{uv} - \boldsymbol{\mu})^\top, \quad (9)$$

where $\boldsymbol{\mu}$ is the mean of all feature vectors, and $\mathbf{X} \in \mathcal{S}_{++}^5$ is an SPD matrix. A set of covariance descriptors $\{\mathbf{X}_i \in \mathcal{S}_{++}^5\}_{i=1}^n$ is the training data, and centers are represented by $\{M_j \in \mathcal{S}_{++}^5\}_{j=1}^m$. In the task, the prediction function $f(\cdot)$ is the affine invariant metric (AIM) [4], which computes distances between centers and descriptors. The distance between a center M_j and a descriptor \mathbf{X}_i is

$$f(M_j, \mathbf{X}_i) = d(M_j, \mathbf{X}_i) = \left\| \logm(\mathbf{X}_i^{-\frac{1}{2}} M_j \mathbf{X}_i^{-\frac{1}{2}}) \right\|_F, \quad (10)$$

where $\|\cdot\|_F$ is the Frobenius-norm, and $\logm(\cdot)$ is the matrix logarithmic function. The objective function $l(\cdot)$ is $l(f(M_j, \mathbf{X}_i)) = f(M_j, \mathbf{X}_i)^2$. The optimizer discovers the centers via minimizing distances between each descriptor and the center that the descriptor belongs to, given by

$$\begin{aligned} \min_{\{M_j \in \mathcal{S}_{++}^d\}} \mathcal{L}(\{M_j\}) &\triangleq \frac{1}{n} \sum_i^n d(M_{c(i)}, \mathbf{X}_i)^2 \\ &= \frac{1}{n} \sum_i^n \left\| \logm(\mathbf{X}_i^{-\frac{1}{2}} M_{c(i)} \mathbf{X}_i^{-\frac{1}{2}}) \right\|_F^2, \end{aligned} \quad (11)$$

where $\mathbf{M}_{c^{(i)}}$ denotes the center that \mathbf{X}_i belongs to. Since we need to solve multiple centers in the clustering tasks, we do not consider to optimize multiple SPD parameters in the inner loop, and the number of centers is regarded as the batchsize of training our optimizer.

The meta-objective to learn our optimizer is

$$\begin{aligned} \min_{\phi} \mathcal{J}(\phi) &= \sum_t^T \mathcal{L}(\{\mathbf{M}_j^{(t+1)}\}) \\ &= \frac{1}{n} \sum_{t,i}^{T,n} \left\| \logm(\mathbf{X}_i^{-\frac{1}{2}} \mathbf{M}_{c^{(i)}}^{(t+1)} \mathbf{X}_i^{-\frac{1}{2}}) \right\|_F^2 \\ &= \frac{1}{n} \sum_{t,i}^{T,n} \left\| \logm\left(\mathbf{X}_i^{-\frac{1}{2}} \left(\Gamma_{\mathbf{M}_{c^{(i)}}^{(t)}}(-g_{\phi}(\nabla_{\mathbf{M}_{c^{(i)}}^{(t)}}, \mathbf{S}_j^{(t-1)}))\right) \mathbf{X}_i^{-\frac{1}{2}}\right) \right\|_F^2. \end{aligned} \quad (12)$$

Minimizing Eq. (12) and obtaining the optimal parameter ϕ for the clustering task are similar to those in the metric nearness task.

1.3. Similarity Learning

Similarity learning aims to learn a Mahalanobis metric, through which similar samples have small distances and dissimilar samples have large distances. In similarity learning experiments, we extract a feature vector $\mathbf{x}_i \in \mathbb{R}^d$ from each image, pairs of vectors $(\mathbf{x}_i, \mathbf{x}_{i'})$ are the training data, and their labels are $y_{ii'}$. $y_{ii'} = 1$ means \mathbf{x}_i and $\mathbf{x}_{i'}$ are similar; otherwise, $y_{ii'} = 0$. In this task, the prediction function $f(\cdot)$ is the Mahalanobis metric that computes the distance between \mathbf{x}_i and $\mathbf{x}_{i'}$ by

$$\begin{aligned} f(\mathbf{M}, \mathbf{x}_i, \mathbf{x}_{i'}) &= d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_{i'}) \\ &= \sqrt{(\mathbf{x}_i - \mathbf{x}_{i'})^{\top} \mathbf{M} (\mathbf{x}_i - \mathbf{x}_{i'})}, \end{aligned} \quad (13)$$

where $\mathbf{M} \in \mathbb{R}^{d \times d}$ is the metric matrix, and the non-negative condition of the metric requires \mathbf{M} to be an SPD matrix, *i.e.*, $\mathbf{M} \in \mathcal{S}_{++}^d$. We use the contrastive loss as the objective function $l(\cdot)$, given by

$$\begin{aligned} &l(f(\mathbf{M}, \mathbf{x}_i, \mathbf{x}_{i'}), y_{ii'}) \\ &= y_{ii'} \cdot \max\left(f(\mathbf{M}, \mathbf{x}_i, \mathbf{x}_{i'}) - \zeta_s, 0\right)^2 \\ &+ (1 - y_{ii'}) \cdot \max\left(\zeta_d - f(\mathbf{M}, \mathbf{x}_i, \mathbf{x}_{i'}), 0\right)^2. \end{aligned} \quad (14)$$

We expect that the distance between two similar samples is smaller than a threshold ζ_s , and the distance between two dissimilar samples is larger than a threshold ζ_d . The loss

function of the similarity learning task is

$$\begin{aligned} \min_{\mathbf{M} \in \mathcal{S}_{++}^d} \mathcal{L}(\mathbf{M}) &= \frac{1}{n} \sum_{i=1}^n l(f(\mathbf{M}, \mathbf{x}_i, \mathbf{x}_{i'}), y_{ii'}) \\ &= \frac{1}{|\mathcal{S}|} \sum_{i,i' \in \mathcal{S}} y_{ii'} \cdot \max\left(d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_{i'}) - \zeta_s, 0\right)^2 \\ &+ \frac{1}{|\mathcal{D}|} \sum_{i,i' \in \mathcal{D}} (1 - y_{ii'}) \cdot \max\left(\zeta_d - d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_{i'}), 0\right)^2, \end{aligned} \quad (15)$$

where n is the number of pairs, \mathcal{S} is the similar pair set, \mathcal{D} is the dissimilar pair set, and $|\mathcal{S}|$ and $|\mathcal{D}|$ are numbers of pairs in \mathcal{S} and \mathcal{D} , respectively.

We train our optimizer according to the following meta-objective

$$\begin{aligned} \min_{\phi} \mathcal{J}(\phi) &= \frac{1}{m} \sum_t^T \sum_j^m \mathcal{L}(\mathbf{M}_j^{(t+1)}) \\ &= \frac{1}{m} \sum_{t,j}^{T,m} \left(\frac{1}{|\mathcal{S}|} \sum_{i,i' \in \mathcal{S}} y_{ii'} \cdot \max\left(d_{\mathbf{M}_j^{(t+1)}}(\mathbf{x}_i, \mathbf{x}_{i'}) - \zeta_s, 0\right)^2 \right. \\ &+ \left. \frac{1}{|\mathcal{D}|} \sum_{i,i' \in \mathcal{D}} (1 - y_{ii'}) \cdot \max\left(\zeta_d - d_{\mathbf{M}_j^{(t+1)}}(\mathbf{x}_i, \mathbf{x}_{i'}), 0\right)^2 \right) \\ &= \frac{1}{m} \sum_{t,j}^{T,m} \left(\frac{1}{|\mathcal{S}|} \sum_{i,i' \in \mathcal{S}} y_{ii'} \right. \\ &\cdot \max\left(\sqrt{(\mathbf{x}_i - \mathbf{x}_{i'})^{\top} \Gamma_{\mathbf{M}_j^{(t)}}(-g_{\phi}(\nabla_{\mathbf{M}_j^{(t)}}, \mathbf{S}_j^{(t-1)}))(\mathbf{x}_i - \mathbf{x}_{i'})} \right. \\ &- \left. \zeta_s, 0\right)^2 \\ &+ \left. \frac{1}{|\mathcal{D}|} \sum_{i,i' \in \mathcal{D}} (1 - y_{ii'}) \cdot \max\left(\zeta_d - \sqrt{(\mathbf{x}_i - \mathbf{x}_{i'})^{\top} \Gamma_{\mathbf{M}_j^{(t)}}(-g_{\phi}(\nabla_{\mathbf{M}_j^{(t)}}, \mathbf{S}_j^{(t-1)}))(\mathbf{x}_i - \mathbf{x}_{i'})}, 0\right)^2 \right), \end{aligned} \quad (16)$$

Minimizing Eq. (16) and obtaining the optimal parameter ϕ of our optimizer for the similarity learning task are similar to those in the metric nearness task.

References

- [1] Justin Brickell, Inderjit S Dhillon, Suvrit Sra, and Joel A Tropp. The metric nearness problem. *SIAM Journal on Matrix Analysis and Applications*, 30(1):375–396, 2008.
- [2] 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 International Conference on Machine Learning (ICML)*, pages 720–729, 2015.
- [3] Peihua Li, Jiangtao Xie, Qilong Wang, and Wangmeng Zuo. Is second-order information helpful for large-scale visual recognition? In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2070–2078, 2017.

- [4] Xavier Pennec, Pierre Fillard, and Nicholas Ayache. A riemannian framework for tensor computing. *International Journal of Computer Vision (IJCV)*, 66(1):41–66, 2006.
- [5] Lei Wang, Jianjia Zhang, Luping Zhou, Chang Tang, and Wanqing Li. Beyond covariance: Feature representation with nonlinear kernel matrices. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.