
Supplementary Material for Efficient Optimization for Linear Dynamical Systems with Applications to Clustering and Sparse Coding

Wenbing Huang^{1,3}, Mehrtash Harandi², Tong Zhang²

Lijie Fan³, Fuchun Sun³, Junzhou Huang¹

¹ Tencent AI Lab. ;

² Data61, CSIRO and Australian National University, Australia;

³ Department of Computer Science and Technology, Tsinghua University,
Tsinghua National Lab. for Information Science and Technology (TNList);

¹{helendhuang, joehhuang}@tencent.com

²{mehrtash.harandi@data61.csiro.au, tong.zhang@anu.edu.cn}

³{flj14@mails, fcsun@mail}.tsinghua.edu.cn

This supplementary material provides the proofs of Theorems 1, 4 and 5, and presents the full flowcharts of applying PGD (Algorithm 1) for clustering and sparse coding. Moreover, we also introduce more details of the datasets *YUPENN* [1], *DynTex* [2] and *DynTex++* [3] that are applied in our experiments. Finally, we provide additional experimental evaluations to compare the classification accuracies between the projection and Martin kernels.

In this supplementary material, bold capital letters denote matrices (*e.g.*, \mathbf{X}) and bold lower-case letters denote column vectors (*e.g.*, \mathbf{x}). \mathbf{I}_n is the $n \times n$ identity matrix. The orthogonal group is denoted by $\mathcal{O}(n)$, *i.e.*, $\mathcal{O}(n) = \{\mathbf{R} \in \mathbb{R}^{n \times n} | \mathbf{R}\mathbf{R}^T = \mathbf{R}^T\mathbf{R} = \mathbf{I}_n\}$. $\|\cdot\|_1$ is the ℓ_1 norm of a vector; $\|\cdot\|_F$ is the Frobenius norm of a matrix. \mathbf{X}^T denotes the matrix transposition. \mathbf{X}^\vee returns the vectorized elements from the columns of \mathbf{X} . \otimes denotes the Kronecker-product. $d\mathbf{X}$ performs the differential operation on \mathbf{X} .

1 Proofs

Theorem 1. *For any given LDS, the system tuple $(\mathbf{A}, \mathbf{C}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{m \times n}$ and all its equivalent representations have the canonical form $(\mathbf{\Lambda}\mathbf{V}, \mathbf{U})$, where $\mathbf{U} \in \text{ST}(m, n)$, $\mathbf{V} \in \mathcal{O}(n)$ and $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$ is diagonal with the diagonal elements arranged in a descend order, *i.e.* $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$.*

Proof. Let the SVD of \mathbf{A} be $\mathbf{U}_A \mathbf{\Lambda} \mathbf{V}_A^T$. According to P-equivalence (Eq.(3)), we obtain

$$\begin{aligned} (\mathbf{A}, \mathbf{C}) &\sim (\mathbf{U}_A^T \mathbf{A} \mathbf{U}_A, \mathbf{C} \mathbf{U}_A) \\ &\sim (\mathbf{\Lambda} \mathbf{V}_A^T \mathbf{U}_A, \mathbf{C} \mathbf{U}_A) \\ &= (\mathbf{\Lambda} \mathbf{V}, \mathbf{U}), \end{aligned} \tag{21}$$

where $\mathbf{V} = \mathbf{V}_A^T \mathbf{U}_A$ and $\mathbf{U} = \mathbf{C} \mathbf{U}_A$ are orthogonal. □

Several matrix computation properties are applied for the proof of Theorem 4:

1. $d(\mathbf{Y}\mathbf{Z}) := (\mathbf{I} \otimes \mathbf{Y})d\mathbf{Z} + (\mathbf{Z}^T \otimes \mathbf{I})d\mathbf{Y}$;
2. $(\mathbf{A}\mathbf{B}\mathbf{C}) := (\mathbf{C}^T \otimes \mathbf{A})\mathbf{B}$;

3. $(ABC) :^T = B :^T (C \otimes A^T)$
4. $(A \otimes B)^T = (A^T \otimes B^T)$

Interested readers can find more details in <http://www.ee.ic.ac.uk/hp/staff/dmb/matrix/calculus.html>.

For better readability, we repeat Theorem 4 before its proof.

Theorem 4. *Let the extended observability matrices of two LDSs (A_1, C_1) and (A_2, C_2) be O_1 and O_2 , respectively. Furthermore, let $G_{12} = O_1^T O_2 = \sum_{t=0}^{\infty} (A_1^T)^t C_1^T C_2 A_2^t$ be the product-matrix between O_1 and O_2 . Given the gradient of the objective function with respect to the product-matrix $\frac{\partial \Gamma}{\partial G_{12}} \doteq H$, the gradients with respect to the system parameters are*

$$\begin{aligned} \frac{\partial \Gamma}{\partial A_1} &= G_{12} A_2 R_{12}^T, & \frac{\partial \Gamma}{\partial C_1} &= C_2 R_{12}^T, \\ \frac{\partial \Gamma}{\partial A_2} &= G_{12}^T A_1 R_{12}, & \frac{\partial \Gamma}{\partial C_2} &= C_1 R_{12}, \end{aligned} \quad (22)$$

where R_{12} is obtained by solving the following DLE

$$A_1 R_{12} A_2^T - R_{12} + H = 0. \quad (23)$$

Proof. The definition $G_{12} = O_1^T O_2$ implies that

$$A_1^T G_{12} A_2 - G_{12} = -C_1^T C_2. \quad (24)$$

By vectorizing and computing the differential on both sides of Eq. (24), we arrive at

$$\begin{aligned} & A_2^T G_{12}^T \otimes I_n dA_1^T : + I_n \otimes A_1^T G_{12} dA_2 : + C_2^T \otimes I_n dC_1^T : + I_n \otimes C_1^T dC_2 : \\ &= (I_{n^2} - A_2^T \otimes A_1^T) dG_{12} : \end{aligned} \quad (25)$$

Thus,

$$\frac{\partial G_{12} :}{\partial A_1^T :} = (I_{n^2} - A_2^T \otimes A_1^T)^{-1} (A_2^T G_{12}^T \otimes I_n), \quad (26)$$

$$\frac{\partial G_{12} :}{\partial A_2 :} = (I_{n^2} - A_2^T \otimes A_1^T)^{-1} (I_n \otimes A_1^T G_{12}), \quad (27)$$

$$\frac{\partial G_{12} :}{\partial C_1^T :} = (I_{n^2} - A_2^T \otimes A_1^T)^{-1} (C_2^T \otimes I_n), \quad (28)$$

$$\frac{\partial G_{12} :}{\partial C_2 :} = (I_{n^2} - A_2^T \otimes A_1^T)^{-1} (I_n \otimes C_1^T), \quad (29)$$

where the invertibility of the term $(I_{n^2} - A_2^T \otimes A_1^T)$ is guaranteed by the stability of A_1 and A_2 .

Under the vectorized form, we have $\frac{\partial \Gamma}{\partial G_{12} :} = (G_{12} :)^T$. Applying the chain rule, we obtain

$$\frac{\partial \Gamma}{\partial A_1^T :} = \frac{\partial \Gamma}{\partial G_{12} :} \frac{\partial G_{12} :}{\partial A_1^T :} = (R_{12} :)^T (A_2^T G_{12}^T \otimes I_n) = (R_{12} A_2^T G_{12}^T :)^T, \quad (30)$$

$$\frac{\partial \Gamma}{\partial A_2 :} = \frac{\partial \Gamma}{\partial G_{12} :} \frac{\partial G_{12} :}{\partial A_2 :} = (R_{12} :)^T (I_n \otimes A_1^T G_{12}) = (G_{12}^T A_1 R_{12} :)^T, \quad (31)$$

$$\frac{\partial \Gamma}{\partial C_1^T :} = \frac{\partial \Gamma}{\partial G_{12} :} \frac{\partial G_{12} :}{\partial C_1^T :} = (R_{12} :)^T (C_2^T \otimes I_n) = (R_{12} C_2^T :)^T, \quad (32)$$

$$\frac{\partial \Gamma}{\partial C_2 :} = \frac{\partial \Gamma}{\partial G_{12} :} \frac{\partial G_{12} :}{\partial C_2 :} = (R_{12} :)^T (I_n \otimes C_1^T) = (C_1 R_{12} :)^T, \quad (33)$$

where we have defined

$$(R_{12} :)^T = (H :)^T (I_{n^2} - A_2^T \otimes A_1^T)^{-1}. \quad (34)$$

Then,

$$\begin{aligned}
& (\mathbf{R}_{12} :)^T = (\mathbf{H} :)^T (\mathbf{I}_{n^2} - \mathbf{A}_2^T \otimes \mathbf{A}_1^T)^{-1} \\
\Rightarrow & (\mathbf{R}_{12} :)^T (\mathbf{I}_{n^2} - \mathbf{A}_2^T \otimes \mathbf{A}_1^T) = (\mathbf{H} :)^T \\
\Rightarrow & (\mathbf{I}_{n^2} - \mathbf{A}_2 \otimes \mathbf{A}_1)(\mathbf{R}_{12} :) = \mathbf{H} : \\
\Rightarrow & \mathbf{R}_{12} : - (\mathbf{A}_2 \otimes \mathbf{A}_1) \mathbf{R}_{12} := \mathbf{H} : \\
\Rightarrow & \mathbf{R}_{12} : - (\mathbf{A}_1 \mathbf{R}_{12} \mathbf{A}_2^T) := \mathbf{G}_{12} : \\
\Rightarrow & \mathbf{R}_{12} - \mathbf{A}_1 \mathbf{R}_{12} \mathbf{A}_2^T = \mathbf{H} \\
\Rightarrow & \mathbf{A}_1 \mathbf{R}_{12} \mathbf{A}_2^T - \mathbf{R}_{12} + \mathbf{H} = 0
\end{aligned} \tag{35}$$

Substituting the matrix \mathbf{R}_{12} into Eq. (30-33) concludes the proof.

Note that one can directly derive \mathbf{R}_{12} from Eq. (34). However, it will increase the computational complexity drastically as the inversion of $(\mathbf{I}_{n^2} - \mathbf{A}_2^T \otimes \mathbf{A}_1^T)^{-1}$ leads to $O(n^6)$ flops. We recall that our solution by using the DLE (Eq. (35)) only requires $O(n^3)$ flops. \square

Theorem 5. *The update direction in Eq.(18) is a descent direction.*

Proof. We denote the update in Eq.(18) by d . Then $d = \frac{\varepsilon}{\max(\varepsilon, |\lambda_k - \tau \nabla \lambda_k|)} (\lambda_k - \tau \nabla \lambda_k) - \lambda_k$. To prove d is along a descent direction, we need to prove $d^T \nabla \lambda_k < 0$ for small τ . To be specific, $d^T \nabla \lambda_k = -a\tau(\nabla \lambda_k)^T \nabla \lambda_k - (1-a)\lambda_k^T \nabla \lambda_k$, where $0 < a = \frac{\varepsilon}{\max(\varepsilon, |\lambda_k - \tau \nabla \lambda_k|)} \leq 1$.

Thus,

$$\begin{aligned}
d^T \nabla \lambda_k & \leq -a\tau |\nabla \lambda_k|^2 + (1-a) |\lambda_k^T \nabla \lambda_k|, \\
& \leq -a\tau |\nabla \lambda_k|^2 + (1-a) |\lambda_k| |\nabla \lambda_k|, \quad (\text{Cauchy Schwarz inequality}) \\
& \leq -a\tau |\nabla \lambda_k|^2 + (1-a) \varepsilon |\nabla \lambda_k|, \\
& = a(-\tau |\nabla \lambda_k|^2 + (\frac{1}{a} - 1) \varepsilon |\nabla \lambda_k|), \\
& = a(-\tau |\nabla \lambda_k|^2 + (\max(\varepsilon, |\lambda_k - \tau \nabla \lambda_k|) - \varepsilon) |\nabla \lambda_k|), \\
& \leq a(-\tau |\nabla \lambda_k|^2 + (\max(\varepsilon, |\lambda_k| + \tau |\nabla \lambda_k|) - \varepsilon) |\nabla \lambda_k|), \\
& \leq a(-\tau |\nabla \lambda_k|^2 + \tau |\nabla \lambda_k| |\nabla \lambda_k|), \\
& = 0,
\end{aligned} \tag{36}$$

where $d^T \nabla \lambda_k = 0$ if and only if $|\lambda_k| = \varepsilon$, λ_k and $\nabla \lambda_k$ have the opposite directions. \square

2 Gradients of the kernels with respect to the LDS parameters

For the reader's convenience, we also provide the gradients of the projection kernel (Eq. (5)) and the Martin kernel (Eq. (19)). These gradients are necessary to pass the gradients from the loss back to the LDS parameters.

2.1 Projection kernel

Suppose we are updating the r -th dictionary atom and passing the gradient through the kernel between \mathbf{D}_r and \mathbf{D}_j . Recall that the projection kernel is given by

$$k(\mathbf{D}_r, \mathbf{D}_j) = \text{Tr}(\mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr}).$$

Then,

$$\frac{\partial k(\mathbf{D}_r, \mathbf{D}_j)}{\partial \mathbf{G}_{rr}} = -\mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr} \mathbf{G}_{rr}^{-1}, \tag{37}$$

$$\frac{\partial k(\mathbf{D}_r, \mathbf{D}_j)}{\partial \mathbf{G}_{rj}} = \mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1}, \tag{38}$$

$$\frac{\partial k(\mathbf{D}_r, \mathbf{D}_j)}{\partial \mathbf{G}_{jr}} = \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr} \mathbf{G}_{rr}^{-1}. \tag{39}$$

2.2 Martin kernel

The Martin kernel is defined as

$$k(\mathbf{D}_r, \mathbf{D}_j) = \det(\mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr}).$$

Thus,

$$\frac{\partial k(\mathbf{D}_r, \mathbf{D}_j)}{\partial \mathbf{G}_{rr}} = -\det(\mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr}) \mathbf{G}_{rr}^{-1}, \quad (40)$$

$$\frac{\partial k(\mathbf{D}_r, \mathbf{D}_j)}{\partial \mathbf{G}_{rj}} = \det(\mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr}) \mathbf{G}_{jr}^{-1}, \quad (41)$$

$$\frac{\partial k(\mathbf{D}_r, \mathbf{D}_j)}{\partial \mathbf{G}_{jr}} = \det(\mathbf{G}_{rr}^{-1} \mathbf{G}_{rj} \mathbf{G}_{jj}^{-1} \mathbf{G}_{jr}) \mathbf{G}_{jr}^{-1}. \quad (42)$$

3 Algorithms for clustering and sparse coding

In the paper, § 4 has demonstrated how to apply the PGD method to compute the mean for clustering and learn the dictionary atoms for sparse coding. We now embed the PGD method into the implementations of these two tasks and provide full details in Algorithms 2 and 3 below.

Algorithm 2 The PGD method for clustering

Input: The data tuples $\{(\mathbf{A}_i, \mathbf{C}_i)\}_{i=1}^N$; the initialization of the means $\{(\mathbf{A}_{m_i}, \mathbf{C}_{m_i})\}_i^C$;
 According to Theorem 1, compute the canonical formulations of $\{(\mathbf{A}_i, \mathbf{C}_i)\}_{i=1}^N$ and $\{(\mathbf{A}_{m_i}, \mathbf{C}_{m_i})\}_i^C$ as $\{(\mathbf{\Lambda}_i, \mathbf{V}_i, \mathbf{U}_i)\}_{i=1}^N$ and $\{(\mathbf{\Lambda}_{m_i}, \mathbf{V}_{m_i}, \mathbf{U}_{m_i})\}_i^C$, respectively;
for $t = 1$ **to** maxIter **do**
 Assign the data tuples to the closest clusters according to the given metric;
 for $i = 1$ **to** C **do**
 Update the mean tuple $(\mathbf{\Lambda}_{m_i}, \mathbf{V}_{m_i}, \mathbf{U}_{m_i})$ of the i -th cluster via Algorithm 1;
 end for
end for
Output: the means $\{(\mathbf{A}_{m_i}, \mathbf{C}_{m_i})\}_i^C$.

Algorithm 3 The PGD method for sparse coding

Input: The data tuples $\{(\mathbf{A}_i, \mathbf{C}_i)\}_{i=1}^N$; the initialization of the dictionary atoms $\{(\mathbf{A}'_j, \mathbf{C}'_j)\}_j^J$;
 According to Theorem 1, compute the canonical formulations of $\{(\mathbf{A}_i, \mathbf{C}_i)\}_{i=1}^N$ and $\{(\mathbf{A}_{m_i}, \mathbf{C}_{m_i})\}_i^C$ as $\{(\mathbf{\Lambda}_i, \mathbf{V}_i, \mathbf{U}_i)\}_{i=1}^N$ and $\{(\mathbf{\Lambda}'_j, \mathbf{V}'_j, \mathbf{U}'_j)\}_j^J$, respectively;
for $t = 1$ **to** maxIter **do**
 Compute the sparse codes \mathbf{z}_i given LDS dictionary by the homotopy-LARS algorithm [4];
 for $r = 1$ **to** J **do**
 Update the r -th atom via Algorithm 1 with only one iteration;
 end for
end for
Output: the dictionary atoms $\{(\mathbf{A}'_j, \mathbf{C}'_j)\}_j^J$.

4 Datasets

YUPENN dataset introduced in [1] consists of fourteen dynamic scene categories where each category has 30 color videos. Analysing this dataset is challenging as the videos are obtained from various sources, *e.g.*, YouTube, BBC Motion Gallery and Getty Images. The videos have an average dimension of $250 \times 370 \times 145$ and vary significantly in resolution, frame rate, scene appearance, scale, illumination condition, and camera viewpoint. Representative examples from different classes are illustrated in Figure 1. We convert all videos to gray-scales and down-sample each frame to have a maximum spatial dimension of 128 pixels while keeping the original aspect ratio.

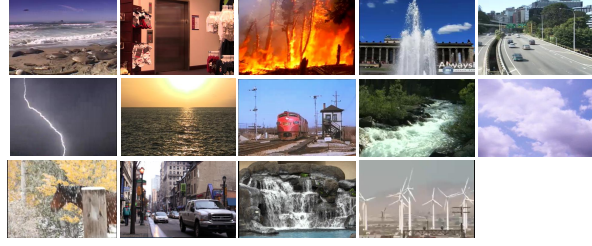


Figure 1: Examples of the YUPENN dataset demonstrating various dynamic scenes (e.g., beach, elevator and forest fire).

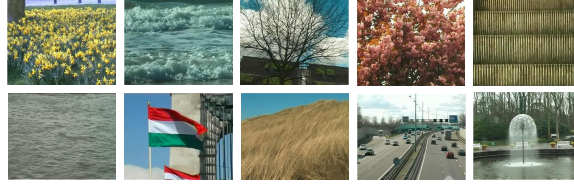


Figure 2: Examples from the DynTex dataset.

The *DynTex* dataset [2] contains $352 \times 288 \times 250$ videos recorded under different environmental conditions, scales and rotations (as illustrated in Figure 2). Three subsets, *i.e.*, *Alpha*, *Beta* and *Gamma* have been applied for classification benchmark in previous studies [2, 5]. However, both *Alpha* and *Beta* have a small number of videos, *i.e.*, 60 and 162, respectively; performing evaluations on them could be easily bias. Hence, we formulate a new dataset by combining the samples of the three subsets, leading to a larger dataset containing 307 videos of 12 classes: *Calm water*, *Escalator*, *Flags*, *Rotation*, *Sea*, *Smoke*, *Traffic*, *Fountain*, *Naked trees*, *Foliage*, *Grass* and *Flowers*. In particular, we first combine the classes *Trees* from *Alpha* and *Beta*, *Naked trees* and *Foliage* from *Gamma* into two non-overlap categories *Naked trees* and *Foliage*; and then combine the classes *Grass* from *Alpha*, *Vegetation* from *Beta*, and *Flowers* and *Grass* from *Gamma* into two categories *Flowers* and *Grass*; and finally integrate other classes of the three subsets. We resize all the videos to 128×128 of gray scale.

DynTex++ [3] is an variant of *DynTex*, where the samples are extracted from the local regions of the videos in *DynTex*. It consists of 3600 videos of 36 classes with 100 videos of size $50 \times 50 \times 50$ per class.

5 Experimental comparison between different kernels

We have performed the classification experiments of PGD based on the projection kernel in the paper (Section 5.2). Now we provide additional experimental evaluations of PGD based on the Martin kernel. In practice, we find that the original definition of Martin kernel in Eq.(19) produces a very small value due to the determinant calculation. We thus revise the Martin kernel by computing powers of the original kernel, namely,

$$k_m((A_1, C_1), (A_2, C_2)) = \left(\det \left(G_{11}^{-1} G_{12} G_{22}^{-1} G_{21} \right) \right)^{\frac{1}{\gamma}}, \quad (43)$$

where γ is set to be 100. We follow the same experimental set-ups as those in the paper. The classification accuracies of the projection kernel and Martin kernel are provided in Table 1. It is observed that Martin kernel almost works better in conjunction with PGD compared with the projection kernel.

Table 1: Mean classification accuracies of the projection and Martin kernels.

Datasets	+BoS		+STPM	
	Projection	Martin	Projection	Martin
YUPENN	84.1	86.1	93.6	94.2
DynTex	65.4	66.7	76.5	74.5

References

- [1] Konstantinos G Derpanis, Matthieu Lecce, Kostas Daniilidis, and Richard P Wildes. Dynamic scene understanding: The role of orientation features in space and time in scene classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1306–1313. IEEE, 2012.
- [2] Renaud Péteri, Sándor Fazekas, and Mark J. Huiskes. DynTex : a Comprehensive Database of Dynamic Textures. *Pattern Recognition Letters*, doi: 10.1016/j.patrec.2010.05.009, 2010. <http://projects.cwi.nl/dyntex/>.
- [3] Bernard Ghanem and Narendra Ahuja. Maximum margin distance learning for dynamic texture recognition. In *European Conference on Computer Vision (ECCV)*, pages 223–236. Springer, 2010.
- [4] David L Donoho and Yaakov Tsaig. Fast solution of ℓ_1 -norm minimization problems when the solution may be sparse. *IEEE Transactions on Information Theory*, 54(11):4789–4812, 2008.
- [5] Yuhui Quan, Chenglong Bao, and Hui Ji. Equiangular kernel dictionary learning with applications to dynamic texture analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 308–316, 2016.