# User Guide

Si-Sheng Young[†]

April 9, 2024

Let us explain how to use the SISHY algorithm for solving the principal component analysis (PCA) problem from incomplete data. Since SISHY has three modes, i.e., "Mode1", "Mode2" **(suggested default mode)** , and "Mode3", that correspond to "Algorithm1", "Algorithm2", and "Algorithm3" in our journal paper, respectively. Note that the closed-form solution of "Algorithm3" is efficiently handled by the so-called plug-and-play (PnP) strategy thereby we employ the pretrained denoising convolutional neural network, "DnCNN" for this demo. However, users can replace "DnCNN" to any denoiser for their own implementation.

## 1 Prerequisites

The SISHY is tested by MATLAB R2022a under Windows 10. Besides, please install the following toolbox for a complete demonstration:

- Deep Learning Toolbox™. (Only for users desiring to implement "Mode3")

- Probabilistic PCA and Factor Analysis.

## 2 Run the code

The "main simulation.m" and "main.m" files demonstrate the subspace identification process on the simulation data and hyperspectral data, respectively.

## 3 Citation

If you find our work useful in your research or publication, please kindly cite our work:

- @ARTICLE10474406,
  author={Lin, Chia-Hsiang and Young, Si-Sheng},
  journal={IEEE Transactions on Geoscience and Remote Sensing},
  title={Signal Subspace Identification for Incomplete Hyperspectral Image With Applications to Various Inverse Problems},
  year={2024},
  volume={62},
  number={},
  pages={1-16},
  doi={10.1109/TGRS.2024.3378705}

---

[†]Institute of Computer and Communication Engineering, National Cheng Kung University, Tainan, Taiwan (R.O.C.) E-mail: `q38121509@gs.ncku.edu.tw`; Website: `https://sites.google.com/view/chiahsianglin/home`

# 4    Mathematical Model

For better user comprehension, we concisely recall the mathematical model of the renowned signal subspace identification (SSID) algorithm, i.e., principal component analysis (PCA), and the proposed SISHY. Consider the data matrix $\boldsymbol{X} = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_L] \in \mathbb{R}^{M \times L}$ with each data point $\boldsymbol{x}_i \in \mathbb{R}^M$. We assume that $\boldsymbol{X}$ is of full row rank (thus, $M \leq L$). Usually, we have $N < M < L$, where $N$ is the model order. The traditional PCA can be expressed as the following optimization problems:

$$\text{PCA}(\boldsymbol{X}) \equiv (\boldsymbol{E}^\star, \boldsymbol{S}^\star, \boldsymbol{d}^\star) := \arg \min_{\boldsymbol{E} \in \mathbb{R}^{M \times N},\ \boldsymbol{S} \in \mathbb{R}^{N \times L},\ \boldsymbol{d} \in \mathbb{R}^M} \frac{1}{2} \|\boldsymbol{X} - (\boldsymbol{E}\boldsymbol{S} + \boldsymbol{d}\boldsymbol{1}_L^T)\|_F^2,$$

where $\boldsymbol{E}^\star$, $\boldsymbol{S}^\star$, and $\boldsymbol{d}^\star$ denote the optimal basis matrix, coefficient matrix, and the data-mean vector, respectively. However, traditional PCA barely handles the incomplete data matrix. To address this, we write the SSID criterion for the given incomplete data $\boldsymbol{X}_{\boldsymbol{\Omega}} \in \mathbb{R}^{M \times L}$ with observed index $\boldsymbol{\Omega} \in \mathbb{R}^{M \times L}$ as

$$\min_{\boldsymbol{E} \in \mathbb{R}^{M \times N},\ \boldsymbol{S} \in \mathbb{R}^{N \times L},\ \boldsymbol{d} \in \mathbb{R}^M} \sum_{(m,\ell) \in \boldsymbol{\Omega}} \frac{1}{2} [x_{m\ell} - (\boldsymbol{E}_{m,:}\boldsymbol{s}_\ell + d_m)]^2. \tag{1}$$

where $\boldsymbol{S} = [\boldsymbol{s}_1, \ldots, \boldsymbol{s}_L]$, $x_{m\ell} \triangleq [\boldsymbol{X}]_{m,\ell}$, $\boldsymbol{E}_{m,:}$ is the $m$th row of $\boldsymbol{E}$, and $d_m$ is the $m$th entry of $\boldsymbol{d}$. Furthermore, we incorporate an auxiliary variable $\boldsymbol{Z} \in \mathbb{R}^{M \times L}$ together with an additional constraint, into the proposed SSID criterion, leading (1) to

$$\text{SISHY}(\boldsymbol{X}_{\boldsymbol{\Omega}}) \equiv (\boldsymbol{E}^*, \boldsymbol{S}^*, \boldsymbol{d}^*, \boldsymbol{Z}^*) := \arg \min_{\boldsymbol{E} \in \mathbb{R}^{M \times N},\ \boldsymbol{S} \in \mathbb{R}^{N \times L},\ \boldsymbol{d} \in \mathbb{R}^M,\ \boldsymbol{Z} \in \Omega_{\boldsymbol{Z}}} \frac{1}{2} \|\boldsymbol{Z} - (\boldsymbol{E}\boldsymbol{S} + \boldsymbol{d}\boldsymbol{1}_L^T)\|_F^2,$$

more details can be found in our journal paper. Accordingly, we summarize the input and output of the proposed SISHY as follows:

- Input

    1. $\boldsymbol{X}_{\boldsymbol{\Omega}} \in \mathbb{R}^{M \times L}$ is the input data matrix with **missing entries set as 0**.
    2. $N$ is the model order/number of principal components.

- Output

    1. $\boldsymbol{E}^* \in \mathbb{R}^{M \times N}$ is the basis matrix whose columns are PCs.
    2. $\boldsymbol{S}^* \in \mathbb{R}^{N \times L}$ coefficient matrix (dimension-reduced data matrix).
    3. $\boldsymbol{d}^* \in \mathbb{R}^M$ is the data-mean vector.
    4. $\boldsymbol{Z}^* \in \mathbb{R}^{M \times L}$ reconstructed data matrix.

# 5    Frequently Asked Question

To enhance users' understanding, we also organize several questions as hereafter:

- **Q:** What representation should we use to denote those missing entries?

    **A:** In our SISHY algorithm, users can simply set the missing entries to 0. The algorithm will consider the 0 as missing entries.

- **Q:** If we input a specific incomplete data matrix, how can we obtain the dimension-reduced data matrix?

  **A:** Our algorithm will automatically output the coefficient matrix $\boldsymbol{S}^*$ corresponding to the input data matrix, which is the dimension-reduced data matrix.

- **Q:** After we process the dimension-reduced data matrix, how can I project the processed dimension-reduced data matrix $\widehat{\boldsymbol{S}}$ to the original dimension?

  **A:** We remark that the basis matrix $\boldsymbol{E}^*$ obtained by SISHY is both full-column-rank and orthonormal (i.e., $\boldsymbol{E}^{*T}\boldsymbol{E}^* = \boldsymbol{I}_N$) thereby we can simply multiply $\widehat{\boldsymbol{S}}$ by $\boldsymbol{E}^*$ followed by plus-ing the data-mean vector $\boldsymbol{d}^*$ to obtain the processed original-dimension data matrix $\widehat{\boldsymbol{X}}$, i.e., $\widehat{\boldsymbol{X}} = \boldsymbol{E}^*\widehat{\boldsymbol{S}} + \boldsymbol{d}^*\mathbf{1}_L^T$, where $\mathbf{1}_L$ is all one vector with length $L$.

- **Q:** If our incomplete data is a 3D tensor (e.g., image), how can we use SISHY to obtain the dimension-reduced data

  **A:** In this case, we need to "reshape" the data tensor into matrix form followed by in-putting the reshaped data matrix into SISHY. We remark that in our "main.m" files, we have demonstrated how to reshape a 3D data cube (i.e., hyperspectral image) into the matrix, so users can simply follow our steps.

- **Q:** What is the auxiliary variable $\boldsymbol{Z}^*$ use for?

  **A:** The auxiliary variable $\boldsymbol{Z}^*$ obtained from SISHY is the reconstructed data matrix. More specifically, if only a few missing entries within the input data matrix, we can employ SISHY to rebuild the complete data matrix. We also remark that in our "main.m" files, we visualize the $\boldsymbol{Z}^*$ indeed is a complete hyperspectral image.