

Minimum-Volume Weighted Symmetric Nonnegative Matrix Factorization for Clustering

Tianxiang Gao[†], Sigurdur Olafsson[†], Songtao Lu[‡]

[†]Department of IMSE, Iowa State University

[‡]Department of ECE, Iowa State University

{gaotx, olafsson, songtao}@iastate.edu



Abstract

In recent years, nonnegative matrix factorization (NMF) attracts much attention in machine learning and signal processing fields due to its interpretability of data in a low dimensional subspace. For clustering problems, symmetric nonnegative matrix factorization (SNMF) as an extension of NMF factorizes the similarity matrix of data points directly and outperforms NMF when dealing with nonlinear data structure. However, the clustering results of SNMF is very sensitive to noisy data. In this paper, we propose a *minimum-volume weighted SNMF* (MV-WSNMF) based on the relationship between robust NMF and SNMF. The proposed MV-WSNMF can approximate the similarity matrices flexibly such that the resulting performance is more robust against noise. A computationally efficient algorithm is also proposed with convergence guarantee. The numerical simulation results show the improvement of the proposed algorithm with respect to clustering accuracy in comparison with the state-of-the-art algorithms.

Introduction

Unsupervised learning problems can be understood as factorizing a data matrix X via factors U and V subject to different constrains

$$\begin{aligned} & \underset{U, V}{\text{minimize}} \quad \|X - UV^T\|_F^2 \\ & \text{subject to} \quad V \in \{0, 1\}^{N \times K}, \mathbf{1}_K^T V = \mathbf{1}_N, \end{aligned} \quad (1)$$

NMF and dimension reduction have the same formulation but with various constrains. For example, NMF relaxes discreteness but enforces nonnegativity for a parts-based representation.

Main Contributions

- A minimum-volume WSNMF method is proposed, which shows strong robustness against noise and provides accurate clustering results.
- An efficient four-block cyclic optimization algorithm is derived for solving the MV-WSNF problem, where every limit point converges to a stationary point.

Minimum-Volume Enclosing Simplex

In *hyperspectral unmixing* (HU) problem, MVES [3] is to circumscribe the data as compact as possible, which is formulated as follows

$$\begin{aligned} & \underset{U, Z}{\text{minimize}} \quad \|X - UZ\|_F^2 + \beta \text{vol}(U) \\ & \text{subject to} \quad Z \geq 0, Z^T \mathbf{1}_K = \mathbf{1}_N, \end{aligned} \quad (2)$$

where U represents a material signature matrix, and $\beta > 0$ is a parameter for taking the tradeoff between the approximation error and volume minimization, and $\text{vol}(U) \triangleq \log \det(U^T U)$ [1].

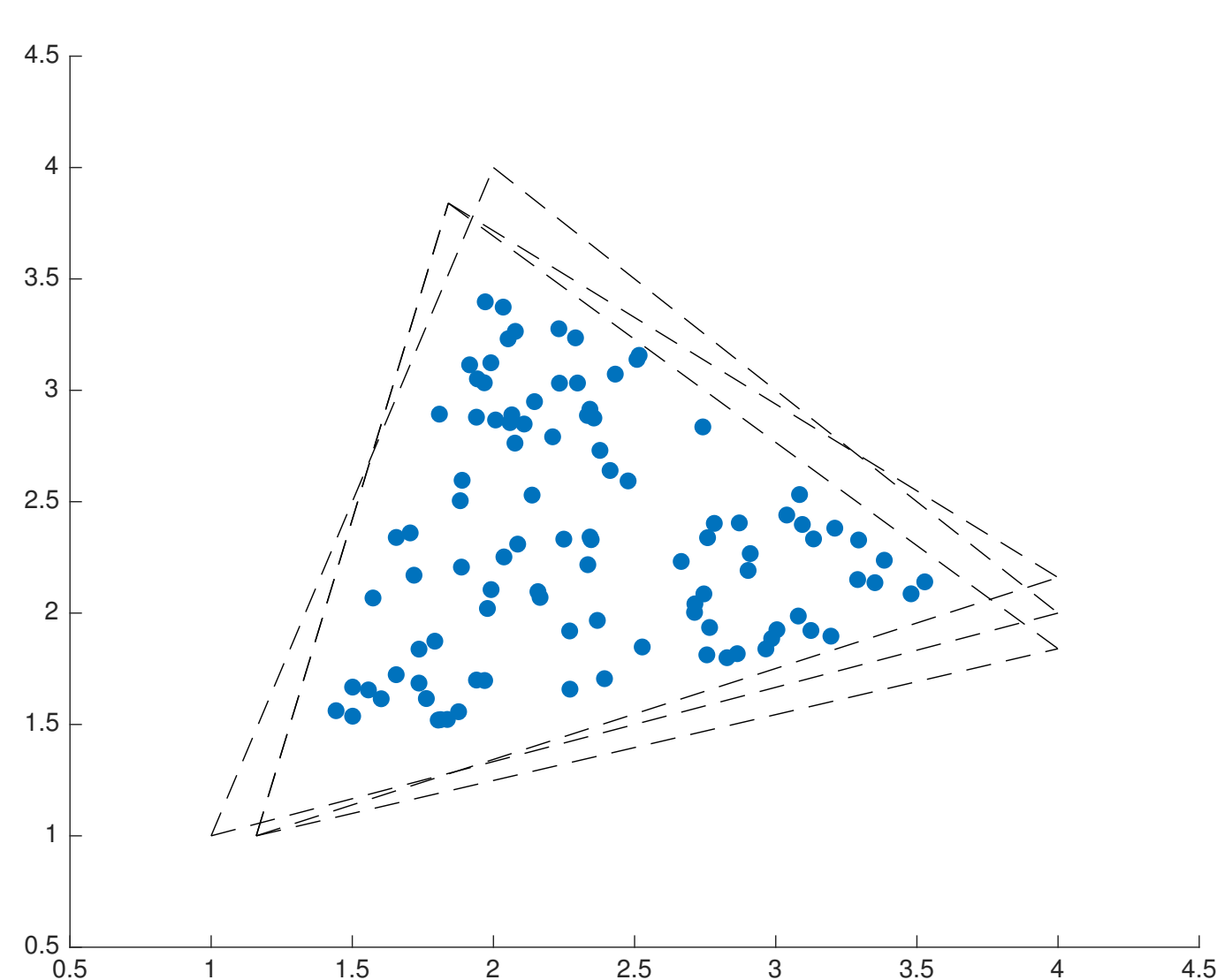


Figure 1: Geometric interpretation of the nonuniqueness of the solution.

Minimum-Volume WSNMF

To conquer the nonlinear data structure with noisy measurements, we propose *minimum-volume WSNMF* (MV-WSNMF) for clustering problems as follows

$$\begin{aligned} & \underset{H, S}{\text{minimize}} \quad \|A - HSH^T\|_F^2 + \beta \log \det(S + \epsilon \mathbf{I}_K) \\ & \text{subject to} \quad H \geq 0, \end{aligned} \quad (3)$$

where A is similarity matrix, and $\epsilon > 0$ is a number in order to prevent the second term unbounded and also ensures $S + \epsilon \mathbf{I}_K$ is *positive definite*.

Variable Splitting for MV-WSNMF

To solve the problem (3), we propose a *variable splitting* algorithm for MV-WSNMF that can converge to a stationary point. The algorithm updates the variables in a cyclic fashion, where the block coordinate descent framework can be directly applied.

A Nonsymmetric Formulation for MV-WSNMF

With leveraging the lemma introduced in [2, Lemma. 2], and splitting the variable H as $H \approx W$, problem (3) can be reformulated as

$$\begin{aligned} & \underset{W, H, S, F}{\text{minimize}} \quad \mathcal{L} = \|A - WSH^T\|_F^2 + \alpha \|W - H\|_F^2 \\ & \quad \quad \quad + \beta (\text{tr}(F(S + \epsilon \mathbf{I}_K)) - \log \det(F)) \\ & \text{subject to} \quad H, W \geq 0, F \geq 0, \end{aligned} \quad (4)$$

where α is the regularization parameter that balances the accurate reconstruction and the difference between W and H .

Algorithm for Nonsymmetric MV-WSNMF

Note that we can cyclically update S , F , H and W until a stopping criteria is satisfied.

Algorithm 1: The Variable Splitting Algorithm

Input: similarity matrix A , number of clusters K

Output: H

initialize H

$W \leftarrow H$

repeat

$$S \leftarrow \frac{1}{2} (W^T W)^\dagger (W^T A H + H^T A W - \beta F) (H^T H)^\dagger$$

$$F \leftarrow (S + \epsilon \mathbf{I}_K)^{-1}$$

$$W \leftarrow \underset{W \geq 0}{\text{argmin}} \left\| \begin{bmatrix} H S \\ \sqrt{\alpha} \mathbf{I}_K \end{bmatrix} W^T - \begin{bmatrix} A \\ \sqrt{\alpha} H^T \end{bmatrix} \right\|_F^2$$

$$H \leftarrow \underset{H \geq 0}{\text{argmin}} \left\| \begin{bmatrix} W S \\ \sqrt{\alpha} \mathbf{I}_K \end{bmatrix} H^T - \begin{bmatrix} A \\ \sqrt{\alpha} W^T \end{bmatrix} \right\|_F^2$$

until stopping criteria satisfied;

Proposition. Every limit point of the solution sequence generated by the Algorithm 1 is a stationary point of problem (4).

Numerical Results

We evaluate the performances of MV-WSNMF on a number of synthetic data sets and compare with the existing state-of-the-art clustering algorithms, including K -Means, NMF, *graph-regularized NMF* (GNMF), *locally consistency concept factorization* (LCCF), normalized Cut (NCut), SNMF, WSNMF, and MV-WSNMF.

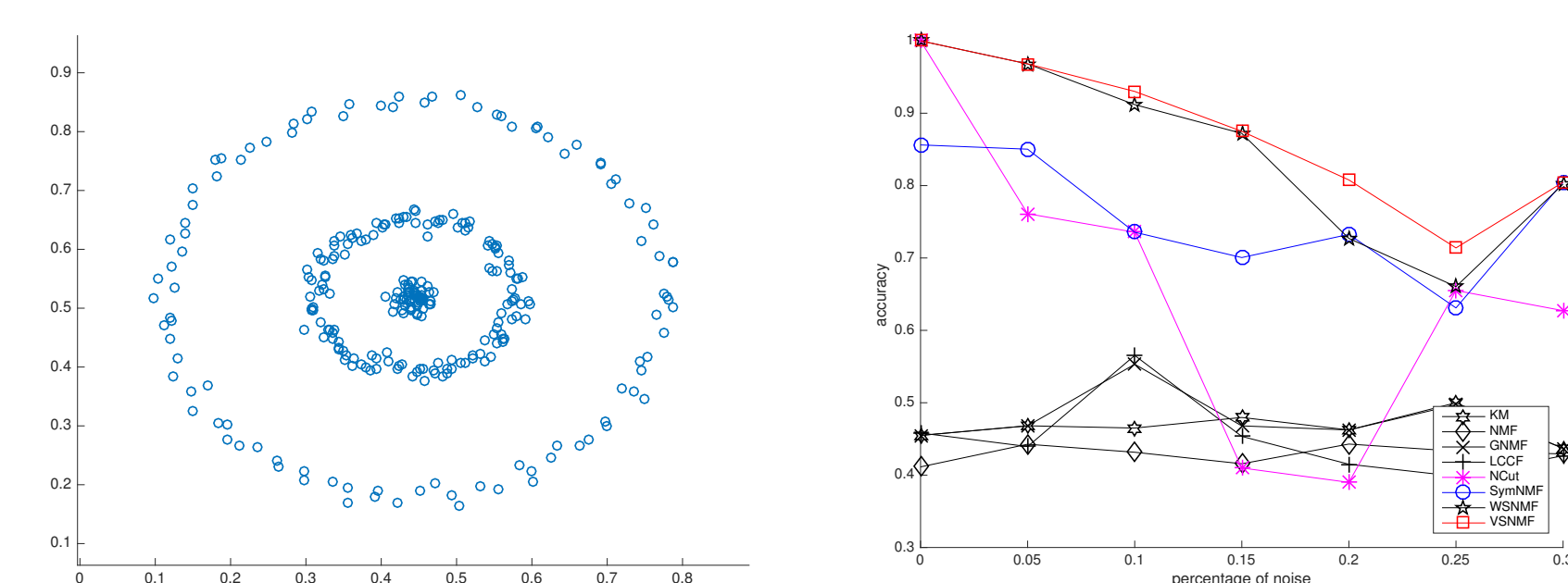


Figure 2: Clustering results for data set 1

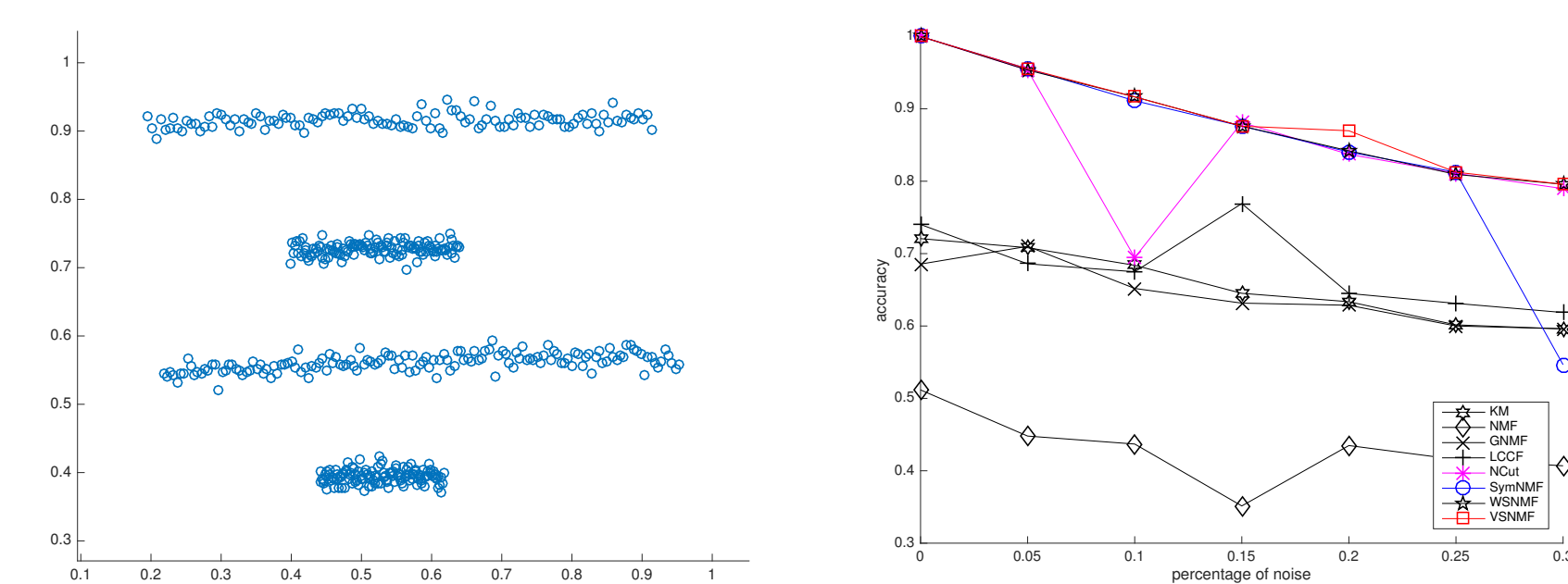


Figure 3: Clustering results for data set 2

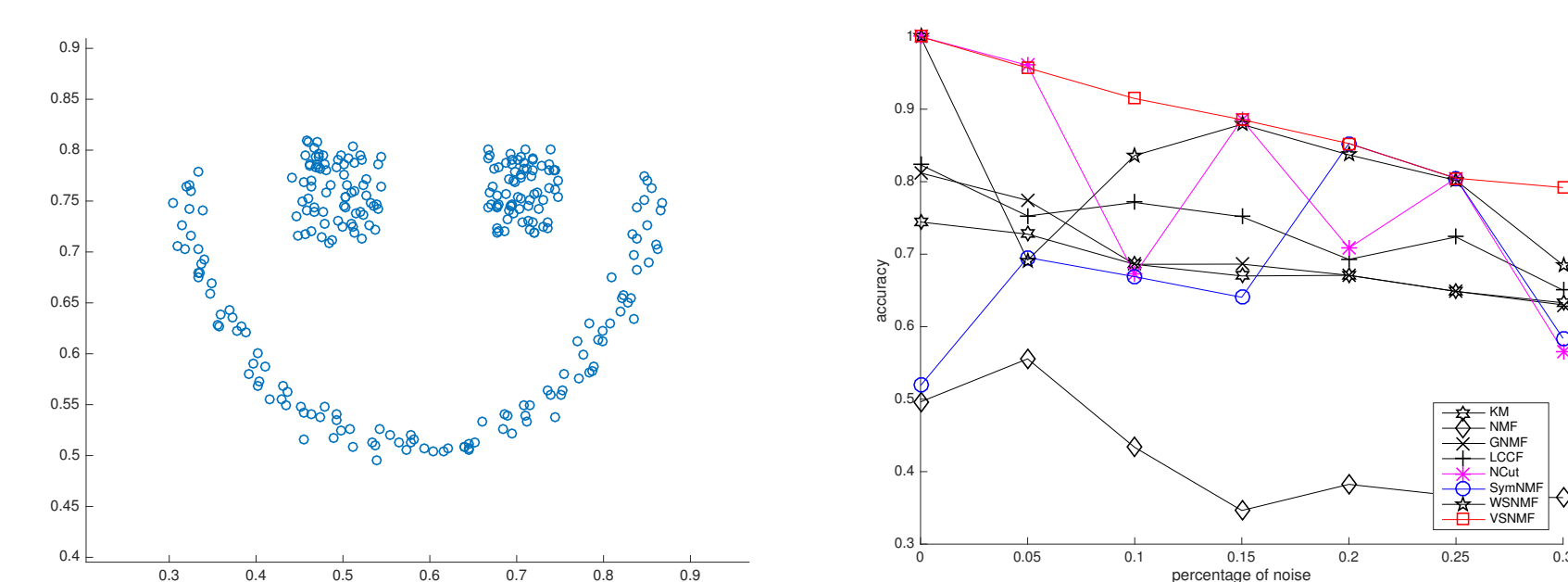


Figure 4: Clustering results for data set 3

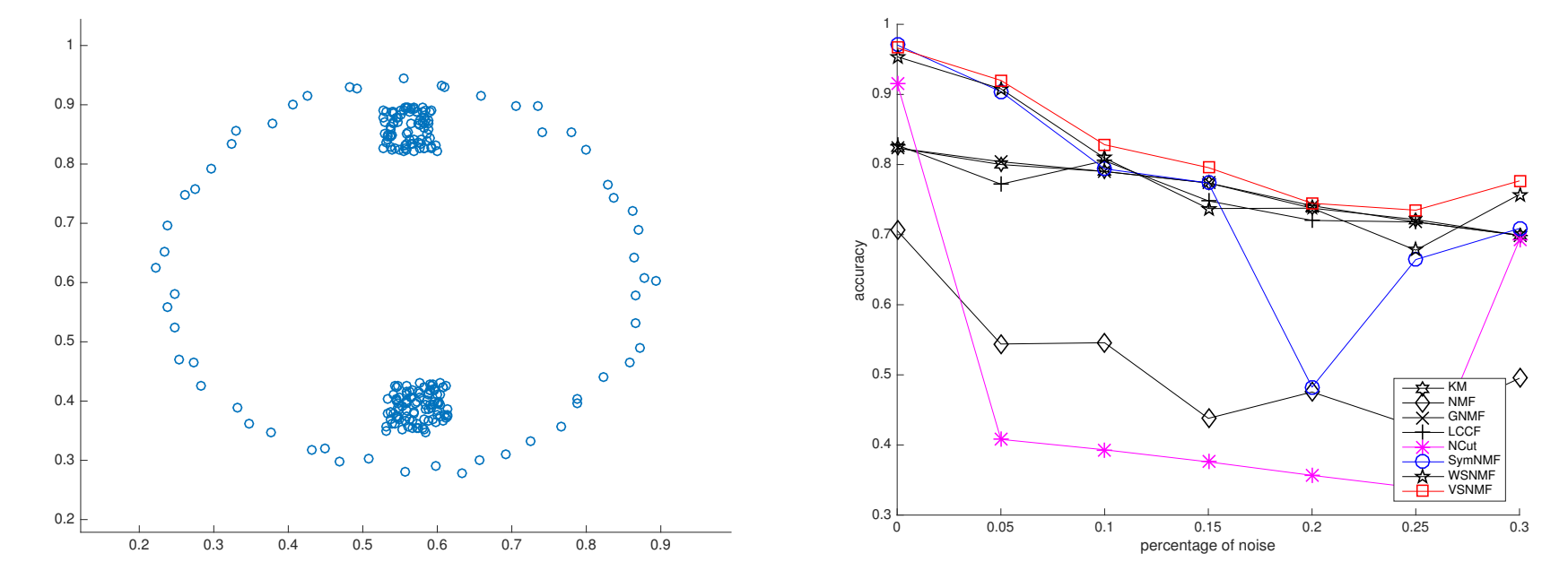


Figure 5: Clustering results for data set 4

Clustering Results

- Kernel functions play an important role in clustering problems since the data sets are transformed in a high dimensional space.
- Factorizing the similarity matrix shows more accurate clustering results. For example, NCut, and SNMF achieve better performance than the NMF-based methods.
- The weighted SNMF including WSNMF and MV-WSNMF usually have better performance in terms of accuracy with or without noise background.
- Regardless of the data sets, the MV-WSNMF achieves the highest accuracy in any percentage of noise data, which shows that MV-WSNMF is stable in noisy data by minimizing the volume of latent space.

It is concluded that with leveraging similarity matrix factorization and minimum volume regularization, MV-SNMF can learn an accurate and compact low dimensional representation of data.

Tuning Parameter Selection

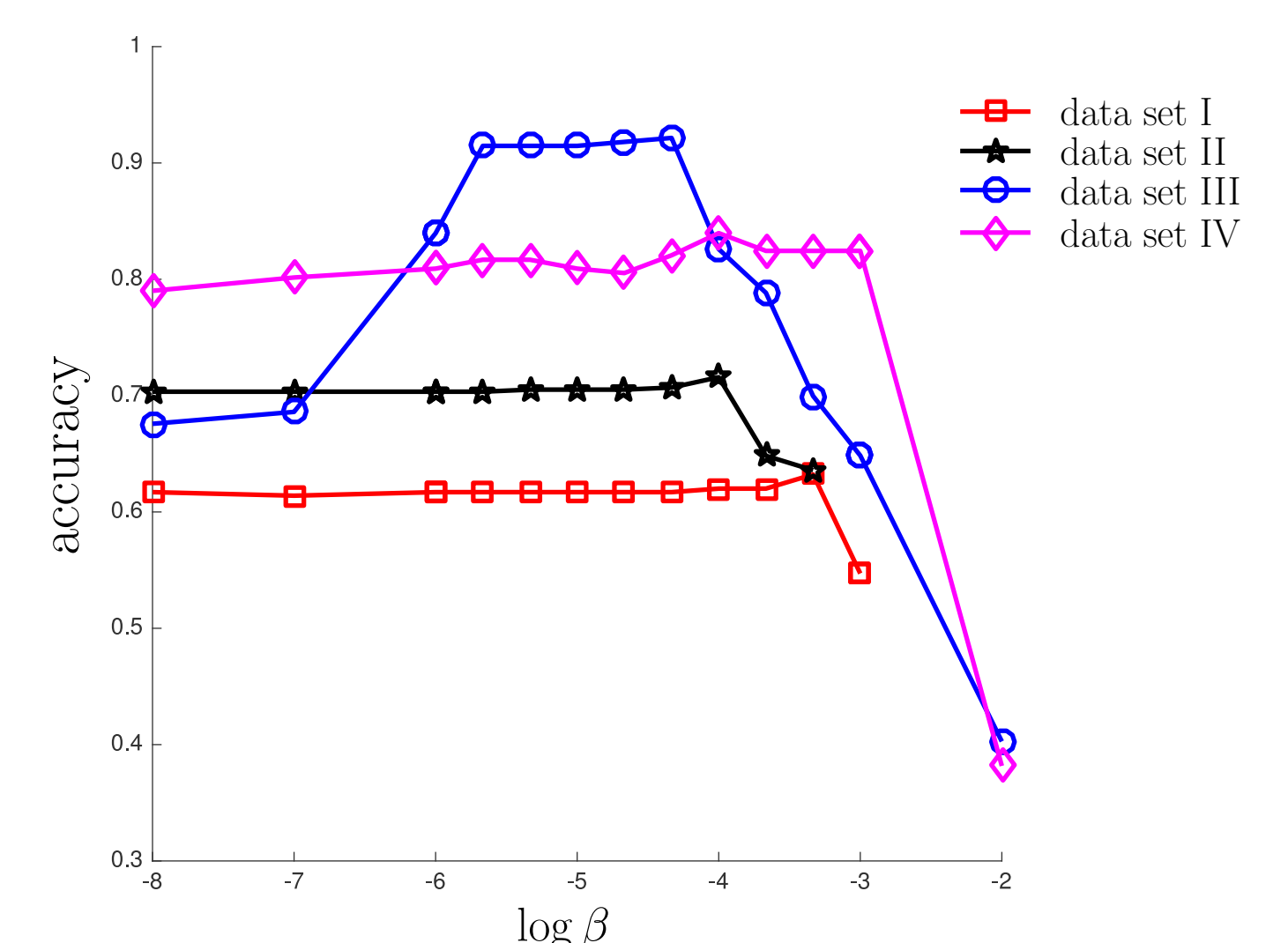


Figure 6: Performance of MV-SNMF varies as the change of β when the percentage of noise data points is 10%

In the MV-SNMF algorithm, there is a regularization parameter β . Figure 6 shows how the clustering accuracy of MV-SNMF changes under various β when the percentage of noise is 10%.

Conclusions

Motivated by the limitation of SNMF, we proposed the MV-WSNMF for improving the robustness of SNMF in clustering problems. An computational efficient four-block cyclical algorithm is proposed and guaranteed to converge to a stationary point. Numerical results show that the clustering accuracy of MV-SNMF based algorithm is much more robust against noisy data in comparison with other benchmark algorithms.

References

- [1] Xiao Fu, Wing-Kin Ma, Kejun Huang, and Nicholas D Sidiropoulos. Robust volume minimization-based matrix factorization via alternating optimization. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2534–2538. IEEE, 2016.
- [2] Jubin Jose, Narayan Prasad, Mohammad Khojastepour, and Sampath Rangarajan. On robust weighted-sum rate maximization in MIMO interference networks. In *Proc. of 2011 IEEE International Conference on Communications (ICC)*, pages 1–6. IEEE, 2011.
- [3] Lidan Miao and Hairong Qi. Endmember extraction from highly mixed data using minimum volume constrained non-negative matrix factorization. *IEEE Transactions on Geoscience and Remote Sensing*, 45(3):765–777, 2007.