MULTI-LINEAR SUBSPACE DISTANCE: A NEW CRITE-RION FOR TENSOR FEATURE SELECTION

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

018

019

020

021

022

024

025

026

027

028

029

031

032

034

037

038

040

041

042 043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Feature selection in tensor data poses greater challenges than in vector representations, since it must capture correlations spanning multiple modes rather than treating each mode in isolation. Existing tensor-based methods partially address this but often treat the feature space as a whole, selecting features globally without respecting mode-specific dependencies. This not only overlooks cross-mode interactions but also increases computational burden, as all features must be considered at once. Moreover, they lack a principled criterion for preserving the global structure of the original tensor. In this work, we introduce Multi-Linear Subspace Learning Feature Selection (MSLFS), a framework that overcomes these limitations by distributing feature selection across modes. Specifically, MSLFS selects a small number of representative slices along each mode, whose intersections yield the most informative features. The core innovation is a multi-linear subspace distance, which provides a principled measure of how well these selected features preserve the global multi-way structure of the data, while significantly reducing redundancy and computational cost. This objective is complemented by two novel regularizations: a joint sparsity constraint that enforces coordinated sparsity across modes to identify compact, non-redundant features, and a higher-order graph constraint that preserves local manifold geometry within the induced subtensor. Taken together, these components guarantee that the overall tensor structure as well as the local neighborhood relationships are preserved. Comprehensive experiments on image recognition and biomedical benchmarks demonstrate that MSLFS consistently surpasses state-of-the-art feature selection techniques in clustering tasks.

1 Introduction

Subspace learning has long served as a foundation for dimensionality reduction, with PCA (Zass & Shashua, 2006), LDA (Jelodar et al., 2019), and their variants (Song et al., 2025; Li et al., 2025) producing low-dimensional embeddings that preserve informative directions. However, these methods operate on vectorized data, discarding multi-way correlations and disrupting the natural geometry of tensorial data such as images and biomedical signals (Liu et al., 2017; Lu et al., 2020). As a result, classical subspace learning often misses key structural dependencies, leading to suboptimal representations for multi-way data (Chouchane et al., 2024).

Recent advances in tensor learning extend linear subspace analysis to multi-way data, enabling models to exploit richer structural information than traditional vector-based methods. Yet, most existing approaches still fall short in how they handle feature selection. In particular, they typically flatten the tensor into a single feature space and select features globally, overlooking the mode-specific dependencies that define the multi-way structure of the data (Chen et al., 2023). This global treatment masks the complementary roles of different modes and forces algorithms to operate over the entire feature set, which becomes computationally expensive in high dimensions. More critically, these methods lack a principled criterion for ensuring that the chosen features preserve the global subspace geometry of the tensor, often capturing only partial correlations.

To overcome these challenges, we introduce *Multi-linear Subspace Learning Feature Selection* (MSLFS), a framework that distributes the selection process across modes rather than treating the feature space as a rigid whole. Instead of picking features globally, MSLFS identifies a small num-

ber of representative slices along each mode; their intersections then form a compact set of features that best reflect the underlying structure of the data. This strategy both respects the multi-way organization of tensor data and reduces computational overhead. At the heart of the framework lies a new notion of *multi-linear subspace distance*, which serves as a principled measure of how well selected features preserve the original multi-way geometry. By optimizing this criterion, MSLFS ensures that the chosen features jointly capture mode-specific information and cross-mode dependencies.

Beyond the core formulation, we introduce two regularizers. The *joint sparsity* term enforces shared sparsity across modes, ensuring that only a compact and representative subset of features is retained. The *higher-order graph* term preserves local manifold geometry in the selected subtensor by extending neighborhood smoothness across all modes. Together, these constraints balance sparsity, global structure, and local geometry. In summary, the contributions of this work are presented as follows.

- A distributed selection strategy is designed to operate across tensor modes, where a small set of representative slices is chosen per mode. Informative features are yielded by their intersections, which respect the multi-way structure while reducing computational cost.
- A novel multi-linear subspace distance is introduced, providing a principled criterion by which the preservation of the global subspace structure across all tensor modes by the selected features is evaluated.
- A joint sparsity constraint is proposed to act simultaneously across multiple tensor modes, whereby a compact and non-redundant subset of features is encouraged while preserving the overall data structure.
- A *higher-order graph* regularization is proposed, through which neighborhood smoothness is extended to tensor data so that local manifold structures are preserved in the reduced representation.

2 RELATED WORK

Vector-Based Unsupervised Feature Selection. Unsupervised feature selection has been widely studied, though most methods target vectorized data rather than multi-dimensional structures. Classical examples include Laplacian Score (LS) (He et al., 2005b), which ranks features by their ability to preserve local geometry. Unsupervised Discriminative Feature Selection (UDFS) (Yang et al., 2011) jointly applies $\ell_{2,1}$ -norm regularization and local discriminative analysis to select sparse and informative features. Sparse PCA for Feature Selection (SPCAFS) (Li et al., 2023) extends PCA with an $\ell_{2,p}$ -norm penalty on the projection matrix, yielding compact feature subsets while retaining principal variance directions.

Tensor-Based Unsupervised Feature Selection. Recently, tensor-based methods have been introduced to overcome the drawbacks of vector-based feature selection, though their use in unsupervised settings remains limited. Among these, two notable approaches have been proposed. Graph Regularized Low-Rank Tensor Representation (GRLTR) (Su et al., 2018) integrates low-rank tensor representation, local geometry preservation, and $\ell_{2,1}$ -norm feature selection, while CPUFS (Chen et al., 2023) combines a tensor-oriented linear classifier, graph-regularized non-negative CP decomposition, and pseudo-label regression. However, these methods still treat the feature space as flat, selecting features globally without considering mode-specific dependencies, which increases computational cost. Our approach instead selects a few representative slices from each mode, whose intersections yield the most informative features. This preserves the multi-way structure, reduces complexity, and ensures the selected features better capture the global data structure.

Notations. For clarity, symbols used in this paper are summarized in Table 1, with detailed descriptions and preliminaries in Appendix 7.1.

3 MULTI-LINEAR SUBSPACE LEARNING

In tensor analysis, multi-linear subspace learning maintains multi-mode structure instead of flattening data (Lu et al., 2011). A major challenge is defining a geometry-aware distance between subspaces spanned by tensor slices across modes. The goal of this section is to define a multi-linear subspace distance which quantifies similarities between these slice-based subspaces, preserv-

Table 1: Summary of notations.

Notation	Meaning
$x, \mathbf{x}, \mathbf{X}, \mathcal{X}$	Scalar; vector; matrix; tensor.
$\mathbf{I}_m, \mathbf{e}_j^{(m)}$	Identity matrix; j -th column.
$\mathbf{A}_{i,:}, \mathbf{\hat{A}}_{:,j}$	i -th row; j -th column of \mathbf{A} .
$\ \mathbf{A}\ _{F}$, $\ \mathbf{A}\ _{2,1}$, $\text{Tr}(\mathbf{A})$	Frobenius norm; $\ell_{2,1}$ -norm; trace.
$\langle \mathbf{u}, \mathbf{v} \rangle, \langle \mathbf{A}, \mathbf{B} \rangle_F$	Dot product; Frobenius inner product.
$\mathbf{A}\odot\mathbf{B},\mathbf{A}\oslash\mathbf{B},\mathbf{A}\otimes\mathbf{B}$	Hadamard product; element-wise division; Kronecker product, where $(\mathbf{A} \otimes \mathbf{B})_{\overline{im}, \overline{jn}} = a_{ij}b_{mn}$.
$\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}, \mathbf{X}_j^{(3)}, \mathbf{X}_{(3)}$	3-mode tensor, with I_3 samples and $I_1 \times I_2$ features; j -th frontal slice; mode-3 unfolding.
$\mathcal{X} \times_n \mathbf{A}, \mathcal{X} \bar{\times}_n \mathbf{v}$	n-mode tensor-matrix; tensor-vector products.
$\operatorname{Ind}^{I_1 \times I_2}$; \mathbb{R}_+	Indicator matrix; Set of non-negative real numbers.

ing cross-mode dependencies and discriminative information. To this end, we first establish the formal definition of the subspace spanned by tensor slices, which serves as the basis for a similarity measure that precisely captures the underlying multi-linear relationships.

Definition 1. Let $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ be a 3-mode tensor with the mode-n slices $\mathbf{X}_1^{(n)}, \dots, \mathbf{X}_{I_n}^{(n)}$, where $n \in \{1, 2, 3\}$. The space spanned by $\mathbf{X}^{(n)} = \{\mathbf{X}_i^{(n)}\}_{i=1}^{I_n}$ is denoted by $\mathcal{S}(\mathbf{X}^{(n)})$ and defined as $\mathcal{S}(\mathbf{X}^{(n)}) = \{\sum_{i=1}^{I_n} \alpha_i^{(n)} \mathbf{X}_i^{(n)} \mid \alpha_i^{(n)} \in \mathbb{R}\}$.

This construction associates each set of tensor slices with a linear subspace, turning the problem of comparing tensor data into a problem of comparing subspaces. To proceed, we need a principled way of measuring how close an external matrix is to such a subspace.

Definition 2. Given $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ and a matrix \mathbf{Z} of the same dimension as a mode-n slice of \mathcal{X} , where $n \in \{1, 2, 3\}$, the distance from \mathbf{Z} to $\mathcal{S}(\mathbf{X}^{(n)})$ is defined as $\mathrm{dist}(\mathbf{Z}, \mathcal{S}(\mathbf{X}^{(n)})) = \min_{\mathbf{W} \in \mathcal{S}(\mathbf{X}^{(n)})} \|\mathbf{Z} - \mathbf{W}\|_F$.

This distance corresponds to the minimum discrepancy between \mathbf{Z} and any element of the subspace. In other words, it quantifies the error incurred when approximating \mathbf{Z} by linear combinations of the mode-n slices of \mathcal{X} . It follows that $\min_{\mathbf{W} \in \mathcal{S}(\mathbf{X}^{(n)})} \|\mathbf{Z} - \mathbf{W}\|_F = \|\mathbf{Z} - \operatorname{Proj}_{\mathcal{S}(\mathbf{X}^{(n)})} \mathbf{Z}\|_F$, where $\operatorname{Proj}_{\mathcal{S}(\mathbf{X}^{(n)})} \mathbf{Z}$ denotes the orthogonal projection of \mathbf{Z} onto the subspace. Since this projection is itself a linear combination of slices, there exists $\alpha^{(n)} = [\alpha_1^{(n)}, \dots, \alpha_{I_n}^{(n)}]^\top \in \mathbb{R}^{I_n}$ such that $\operatorname{Proj}_{\mathcal{S}(\mathbf{X}^{(n)})} \mathbf{Z} = \sum_{i=1}^{I_n} \alpha_i^{(n)} \mathbf{X}_i^{(n)} = \mathcal{X} \bar{\mathbf{x}}_n \alpha^{(n)}$. Consequently, $\operatorname{dist}(\mathbf{Z}, \mathcal{S}(\mathbf{X}^{(n)})) = \|\mathbf{Z} - \mathcal{X} \bar{\mathbf{x}}_n \alpha^{(n)}\|_F$.

Beyond this general case, additional structure yields simplifications. If the slices $\{\mathbf{X}_i^{(n)}\}_{i=1}^{I_n}$ are orthonormal (i.e., $\langle \mathbf{X}_i^{(n)}, \mathbf{X}_j^{(n)} \rangle_F = 0$ for $i \neq j$ and $\|\mathbf{X}_i^{(n)}\|_F = 1$ for all i), the projection coefficients become explicit inner products: $\alpha^{(n)} = [\langle \mathbf{Z}, \mathbf{X}_1^{(n)} \rangle_F, \dots, \langle \mathbf{Z}, \mathbf{X}_{I_n}^{(n)} \rangle_F]^{\top}$. In this case, $\mathrm{dist}(\mathbf{Z}, \mathcal{S}(\mathbf{X}^{(n)})) = \|\mathbf{Z} - \sum_{i=1}^{I_n} \langle \mathbf{Z}, \mathbf{X}_i^{(n)} \rangle_F \mathbf{X}_i^{(n)} \|_F$, which admits a simple geometric interpretation as subtracting the projection of \mathbf{Z} onto the orthonormal basis formed by the mode-n slices.

So far we have defined the distance between a single matrix and the subspace spanned by tensor mode-n slices. Beyond this, the concept can be naturally extended to quantify the distance between two subspaces, each spanned by the mode-n slices of two distinct tensors.

Definition 3 (Multi-linear Subspace Distance). Let $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ be a 3-mode tensor and \mathcal{Y} another 3-mode tensor of the same dimensionality except that its mode-n size equals J_n , where $n \in \{1,2,3\}$. The distance between the mode-n subspaces $\mathcal{S}(\mathbf{X}^{(n)})$ and $\mathcal{S}(\mathbf{Y}^{(n)})$ is defined as $\mathrm{dist}(\mathcal{S}(\mathbf{X}^{(n)}),\mathcal{S}(\mathbf{Y}^{(n)})) = \sum_{i=1}^{I_n} \mathrm{dist}(\mathbf{X}_i^{(n)},\mathcal{S}(\mathbf{Y}^{(n)}))$.

It can be shown that $\operatorname{dist}(\mathcal{S}(\mathbf{X}^{(n)}),\mathcal{S}(\mathbf{Y}^{(n)})) = \sum_{i=1}^{I_n} \|\mathbf{X}_i^{(n)} - \mathcal{Y}\bar{\times}_n\alpha_i^{(n)}\|_F^2 = \|\mathcal{X} - \mathcal{Y}\times_n\mathbf{H}^{(n)}\|_F^2$, where $\mathbf{H}^{(n)} \in \mathbb{R}^{I_n\times J_n}$ is such that its *i*-th row is $\alpha_i^{(n)}$. Thus, the distance admits a compact tensor representation via a reconstruction error term. This formulation essentially measures how far each mode-n slice of \mathcal{X} lies from the $\mathcal{S}(\mathbf{Y}^{(n)})$, and aggregates these deviations across all mode-n slices.

Remark 1. The concept of multi-linear subspace distance provides a key link between tensor geometry and feature selection. Concretely, let $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ denote the tensor data with I_3 samples

 and $I_1 \times I_2$ features. Each mode-3 fiber represents a single feature and can be seen as the intersection of its corresponding mode-1 and mode-2 slices. Thus, the ability of a fiber to characterize the feature space depends on how well these slices span $\mathcal{S}(\mathbf{X}^{(1)})$ and $\mathcal{S}(\mathbf{X}^{(2)})$. The multi-linear subspace distance provides a natural measure to evaluate this, enabling us to identify informative slices across modes whose intersections yield fibers that faithfully preserve the global structure. By minimizing the distance between the full subspace and the one formed by selected slices, our framework ensures fidelity and coherence across modes. This principle provides the foundation for our feature selection strategy, which will be further developed in the following sections.

3.1 SUBTENSORS AND SLICE SELECTION

Building on the idea of multi-linear subspace distance, a natural way to reduce redundancy while preserving structure is to restrict attention to a subset of slices. Such subsets define subtensors, which retain the essential information needed to approximate the span of the full tensor. By working with subtensors, we can formalize slice selection as a principled step in feature selection, preparing the ground for our definition below.

Definition 4. For a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, a subtensor $\mathcal{X}^{(n;k)}$ is obtained by choosing k mode-n slices indexed by $\{i_1^{(n)}, \dots, i_k^{(n)}\}$, where each $i_j^{(n)} \in \{1, \dots, I_n\}$ and $n \in \{1, 2, 3\}$.

Any single mode-n slice $\mathbf{X}_{j}^{(n)}$ can be written as $\mathbf{X}_{j}^{(n)} = \mathcal{X} \bar{\times}_{n} \mathbf{e}_{j}^{(n)}, \forall j \in \{1, \cdots I_{n}\}$, where $\mathbf{e}_{j}^{(n)}$ is the j-th column of the identity $\mathbf{I}_{I_{n}}$. More generally, a subtensor $\mathcal{X}^{(n;k)}$ formed from $\{\mathbf{X}_{i_{1}}^{(n)}, \ldots, \mathbf{X}_{i_{k}}^{(n)}\}$ can be expressed as $\mathcal{X}^{(n;k)} = \mathcal{X} \times_{n} \mathbf{W}^{(n;k)}$, where $\mathbf{W}^{(n;k)} \in \mathbb{R}^{k \times I_{n}}$ is a selection matrix whose rows are standard basis vectors.

Building on this, the distance between the span of all slices and that of a selected subset follows directly. By Definition 3, we obtain

$$\operatorname{dist}(\mathcal{S}(\mathbf{X}^{(n)}), \mathcal{S}(\mathbf{X}^{(n;k)})) = \|\mathcal{X} - \mathcal{X}^{(n;k)} \times_n \mathbf{H}^{(n;k)}\|_F = \|\mathcal{X} - \mathcal{X} \times_n \mathbf{W}^{(n;k)} \times_n \mathbf{H}^{(n;k)}\|_F$$
$$= \|\mathcal{X} - \mathcal{X} \times_n (\mathbf{H}^{(n;k)} \mathbf{W}^{(n;k)})\|_F. \tag{1}$$

This characterization shows that the distances between full and reduced subspaces can be understood as the error of reconstructing the original tensor using only selected slices and suitable weighting.

3.2 Core Representation via Intersection Fibers

The subspace framework developed in (1) can be naturally extended to a compact tensor representation in terms of mode-3 fibers. By selecting slices along modes 1 and 2 that span the corresponding mode subspaces $\mathcal{S}(\mathbf{X}^{(1)})$ and $\mathcal{S}(\mathbf{X}^{(2)})$, we obtain a reduced set of mode-3 fibers located at their intersections. These intersection fibers act as structural representatives, capturing the same subspace as the full collection of mode-3 fibers. Consequently, the entire tensor can be approximated using a core representation derived from this smaller, more informative subset, whose validity is rigorously established by the following theorem.

Theorem 3.1. Let $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ be a 3-mode tensor. Suppose the mode-1 and mode-2 subspaces $\mathcal{S}(\mathbf{X}^{(1)})$ and $\mathcal{S}(\mathbf{X}^{(2)})$ admit bases of dimensions $R_1 \leq I_1$ and $R_2 \leq I_2$, indexed by $T_1 = \{i_1^{(1)}, \dots, i_{R_1}^{(1)}\}$, $T_2 = \{i_1^{(2)}, \dots, i_{R_2}^{(2)}\}$. Let $\mathbf{W}^{(1;R_1)} \in \operatorname{Ind}^{R_1 \times I_1}$ and $\mathbf{W}^{(2;R_2)} \in \operatorname{Ind}^{R_2 \times I_2}$ be the corresponding indicator matrices. For each pair $(i_1, i_2) \in \{1, \dots, I_1\} \times \{1, \dots, I_2\}$, let $\mathbf{f}_{i_1, i_2} = \mathcal{X}_{i_1, i_2, :} \in \mathbb{R}^{I_3}$ denote the mode-3 fiber.

(Part I: Core Dictionary). The R_1R_2 intersection fibers $\{\mathbf{f}_{i_{r_1}^{(1)},i_{r_2}^{(2)}}\}_{r_1,r_2=1}^{R_1,R_2}$ form a *core dictionary* that spans all mode-3 fibers of \mathcal{X} . Stacking them columnwise yields the core matrix

$$\mathbf{F}_{\text{core}} = (\mathcal{X} \times_1 \mathbf{W}^{(1;R_1)} \times_2 \mathbf{W}^{(2;R_2)})_{(3)} = \mathbf{X}_{(3)} (\mathbf{W}^{(2;R_2)} \otimes \mathbf{W}^{(1;R_1)})^{\top}, \tag{2}$$

where $(\mathbf{W}^{(2;R_2)} \otimes \mathbf{W}^{(1;R_1)})^{\top}$ acts as the indicator matrix selecting precisely those core fibers.

(Part II: Separable Reconstruction). There exist coefficient matrices $\mathbf{H}^{(1;R_1)} \in \mathbb{R}^{I_1 \times R_1}$, $\mathbf{H}^{(2;R_2)} \in \mathbb{R}^{I_2 \times R_2}$, such that every mode-3 fiber admits the separable expansion

$$\mathbf{f}_{i_1,i_2} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} h_{i_1,r_1}^{(1;R_1)} h_{i_2,r_2}^{(2;R_2)} \mathbf{f}_{i_{r_1}^{(1)},i_{r_2}^{(2)}},$$
(3)

and equivalently, the unfolding satisfies

$$\mathbf{X}_{(3)} = (\mathcal{X} \times_1 \mathbf{H}^{(1;R_1)} \mathbf{W}^{(1;R_1)} \times_2 \mathbf{H}^{(2;R_2)} \mathbf{W}^{(2;R_2)})_{(3)} = \mathbf{F}_{core}(\mathbf{H}^{(2;R_2)} \otimes \mathbf{H}^{(1;R_1)})^{\top}.$$
(4)

Proof. A detailed proof of this theorem is presented in Appendix 7.2.

Intuition. Fixing bases for $S(\mathbf{X}^{(1)})$ and $S(\mathbf{X}^{(2)})$ encodes the tensor's structure in their R_1R_2 intersection fibers, which act as a compact *core dictionary*, capturing the interactions between the two subspaces. The coefficient matrices $\mathbf{H}^{(1;R_1)}$ and $\mathbf{H}^{(2;R_2)}$ provide separable weights to reconstruct all fibers. Exact recovery is guaranteed when the chosen slices form true bases; otherwise, approximate bases yield reconstructions with errors tied to the residuals.

Remark 2. Theorem 3.1 underpins multi-way feature selection. When modes 1 and 2 correspond to features and mode-3 indexes samples, each mode-3 fiber represents a feature's response across samples. Feature selection thus reduces to choosing representative bases along modes 1 and 2, whose intersection fibers form the most informative representatives of the full feature space.

4 TENSOR-BASED FEATURE SELECTION

In this section, we formalize the task of feature selection in tensor data. The model developed in this section is presented under the assumption that the input is a nonnegative 3-mode tensor. This assumption is well aligned with many practical multi-way datasets such as images, videos, and medical scans, where entries naturally take non-negative values (Bi et al., 2025). Nonetheless, the framework can be readily extended to general tensor data, and we provide a discussion of this extension in Appendix 7.6. Let $\mathcal{X} \in \mathbb{R}_{+}^{I_1 \times I_2 \times I_3}$ be a non-negative 3-mode data tensor with I_3 samples, each described by $I_1 \times I_2$ multiway features. The problem is to select a subset of mode-3 fibers that best preserve the structure of the full tensor.

Feature Selection via Core Theorem. According to Theorem 3.1, this can be achieved by choosing $m_1 \leq I_1$ slices along mode-1 and $m_2 \leq I_2$ slices along

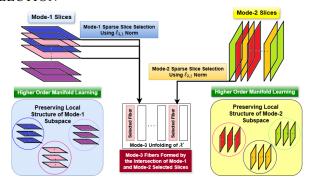


Figure 1: Schematic illustration of multi-linear subspace learning feature selection (MSLFS). Mode-1 and mode-2 slices of the input tensor are processed via $\ell_{2,1}$ -norm based sparse selection, where the joint row sparsity regularization ensures that only a limited number of slice combinations are retained as informative representatives. The first term of the objective function ensures reconstruction fidelity by using the intersection of the selected slices to form representative mode-3 fibers. The second term enforces local manifold preservation within each mode, thereby maintaining the geometric structure of the data subspaces.

mode-2, which approximate $\mathcal{S}(\mathbf{X}^{(1)})$ and $\mathcal{S}(\mathbf{X}^{(2)})$, respectively. The intersection of these selected slices yields a compact yet expressive set of representative mode-3 fibers that best span the feature subspace. Concretely, the feature selection problem can be formulated as follows:

$$\min_{\mathbf{H}^{(1;m_1)},\mathbf{H}^{(2;m_2)},\mathbf{W}^{(1;m_1)},\mathbf{W}^{(2;m_2)} \ge 0} \|\mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2$$
s.t.
$$\mathbf{W}^{(1;m_1)} \in \operatorname{Ind}^{m_1 \times I_1}, \mathbf{W}^{(2;m_2)} \in \operatorname{Ind}^{m_2 \times I_2}.$$
(5)

Here, $\mathbf{W}^{(1;m_1)}$ and $\mathbf{W}^{(2;m_2)}$ are indicator matrices marking the selected slices, and $\mathbf{H}^{(1;m_1)} \in \mathbb{R}^{I_1 \times m_1}$ and $\mathbf{H}^{(2;m_2)} \in \mathbb{R}^{I_2 \times m_2}$ are the corresponding coefficient matrices.

Relaxation via Orthogonality. Since the minimization problem (5) is NP-hard, directly using indicator matrices is impractical. We relax this by enforcing orthogonality on $\mathbf{W}^{(1;m_1)}$ and $\mathbf{W}^{(2;m_2)}$, equivalently on their Kronecker product. Combined with non-negativity, this ensures each column remains one-hot, preserving selection while keeping the optimization tractable.

Row-Sparsity Regularization. Given the sparsity of $\mathbf{W}^{(1;m_1)}$ and $\mathbf{W}^{(2;m_2)}$, their Kronecker product, which acts as the indicator matrix for the $m_1 \times m_2$ intersection mode-3 fibers, inherits this property. To emphasize only the most informative slice combinations, we impose joint row-sparsity on $(\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}$, ensuring that only a few mode-3 fibers dominate the reconstruction and redundancy is reduced. To formalize this idea, we employ the $\ell_{2,1}$ norm. For $(\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}$ this becomes:

$$\|(\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}\|_{2,1} = \text{Tr}((\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)}) \mathbf{U} (\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}), \quad (6)$$

where $\mathbf{U} \in \mathbb{R}^{I_1 I_2 \times I_1 I_2}$ is diagonal with entries equal to the reciprocals of the ℓ_2 norms of the columns of $\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)}$.

Mode-Wise Factorization of the Penalty. Because each column of $\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)}$ is a Kronecker product of a column of $\mathbf{W}^{(2;m_2)}$ and one of $\mathbf{W}^{(1;m_1)}$, the matrix \mathbf{U} decomposes as $\mathbf{U}^{(2)} \otimes \mathbf{U}^{(1)}$, where $\mathbf{U}^{(1)} \in \mathbb{R}^{I_1 \times I_1}$ and $\mathbf{U}^{(2)} \in \mathbb{R}^{I_2 \times I_2}$ are diagonal matrices whose entries depend only on the columns of $\mathbf{W}^{(1;m_1)}$ and $\mathbf{W}^{(2;m_2)}$, respectively. Substituting this gives:

$$\|(\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}\|_{2,1} = \text{Tr}((\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})(\mathbf{U}^{(2)} \otimes \mathbf{U}^{(1)})(\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}).$$
(7)

Using standard Kronecker product identities, this expression simplifies to

$$\text{Tr}((\mathbf{W}^{(2;m_2)}\mathbf{U}^{(2)}\mathbf{W}^{(2;m_2)^{\top}})\otimes (\mathbf{W}^{(1;m_1)}\mathbf{U}^{(1)}\mathbf{W}^{(1;m_1)^{\top}})),$$

and since the trace of a Kronecker product factorizes into the product of traces, we finally obtain

$$\|(\mathbf{W}^{(2;m_2)} \otimes \mathbf{W}^{(1;m_1)})^{\top}\|_{2,1} = \|\mathbf{W}^{(2;m_2)^{\top}}\|_{2,1} \|\mathbf{W}^{(1;m_1)^{\top}}\|_{2,1}$$
$$= \operatorname{Tr}(\mathbf{W}^{(2;m_2)}\mathbf{U}^{(2)}\mathbf{W}^{(2;m_2)^{\top}}) \operatorname{Tr}(\mathbf{W}^{(1;m_1)}\mathbf{U}^{(1)}\mathbf{W}^{(1;m_1)^{\top}}). \quad (8)$$

Interpretation. The $\ell_{2,1}$ penalty factorizes across modes, with each trace term measuring the representational quality of slices in its subspace while penalizing redundancy. This separation reduces computation and enables mode-wise control, ensuring balanced selection that retains only the most informative fibers.

4.1 Graph Regularization for Higher-Order Manifold Learning

In multi-way feature selection, it is crucial to preserve both the global span and the intrinsic geometry of the data. Graph regularization enforces local neighborhood consistency, ensuring proximity in the original space is maintained in the learned representation. Extending this to tensors requires jointly modeling local structures across all modes.

Fiber Representation. Let $\mathbf{H}^{(1;m_1)} \in \mathbb{R}^{I_1 \times m_1}$ and $\mathbf{H}^{(2;m_2)} \in \mathbb{R}^{I_2 \times m_2}$ denote coefficient matrices for the selected slices along modes 1 and 2. By Theorem 3.1, each mode-3 fiber \mathbf{f}_{i_1,i_2} can be approximated in terms of the core fibers as: $\mathbf{f}_{i_1,i_2} = \mathbf{F}_{\text{core}} \left((\mathbf{H}^{(2;m_2)} \otimes \mathbf{H}^{(1;m_1)})^{\top} \right)_{:,\overline{i_1 i_2}}$, where the coefficient vector $\left((\mathbf{H}^{(2;m_2)} \otimes \mathbf{H}^{(1;m_1)})^{\top} \right)_{:,\overline{i_1 i_2}}$ encodes how the fiber is reconstructed from the shared subspace. Intuitively, if two fibers \mathbf{f}_{i_1,i_2} and \mathbf{f}_{j_1,j_2} are similar in the original space, their coefficient vectors should also be close, reflecting their functional similarity in reconstruction.

Graph Regularization. To enforce this locality, we minimize the squared distance between coefficient vectors, weighted by their similarity:

$$\frac{1}{2} \sum_{i_1, i_2} \sum_{j_1, j_2} \left\| \left((\mathbf{H}^{(2; m_2)} \otimes \mathbf{H}^{(1; m_1)})^{\top} \right)_{:, \overline{i_1 i_2}} - \left((\mathbf{H}^{(2; m_2)} \otimes \mathbf{H}^{(1; m_1)})^{\top} \right)_{:, \overline{j_1 j_2}} \right\|_2^2 b_{\overline{i_1 i_2}, \overline{j_1 j_2}}, \quad (9)$$

where $b_{\overline{i_1i_2},\overline{j_1j_2}}$ encodes the similarity between \mathbf{f}_{i_1,i_2} and \mathbf{f}_{j_1,j_2} . This term can be rewritten compactly in matrix form as: $\mathrm{Tr}\big[(\mathbf{H}^{(2;m_2)}\otimes\mathbf{H}^{(1;m_1)})^{\top}\mathbf{L}(\mathbf{H}^{(2;m_2)}\otimes\mathbf{H}^{(1;m_1)})\big]$, where $\mathbf{L}\in\mathbb{R}^{I_1I_2\times I_1I_2}$ is the Laplacian of the feature similarity graph.

Mode-Wise Decomposition. To ease computation, we exploit the fact that similarities between fibers factorize across modes. This induces a Kronecker structure in the joint Laplacian, expressed as $\mathbf{L} = \mathbf{L}^{(2)} \otimes \mathbf{L}^{(1)}$. For each mode $n \in \{1,2\}$, the Laplacian $\mathbf{L}^{(n)} = \mathbf{A}^{(n)} - \mathbf{B}^{(n)}$ is constructed from the degree matrix $\mathbf{A}^{(n)}$ and similarity matrix $\mathbf{B}^{(n)}$, with $\mathbf{L}^{(n)}, \mathbf{A}^{(n)}, \mathbf{B}^{(n)} \in \mathbb{R}^{I_n \times I_n}$. Substituting this decomposition yields:

$$\operatorname{Tr}\left[(\mathbf{H}^{(2;m_2)} \otimes \mathbf{H}^{(1;m_1)})^{\top} (\mathbf{L}^{(2)} \otimes \mathbf{L}^{(1)}) (\mathbf{H}^{(2;m_2)} \otimes \mathbf{H}^{(1;m_1)}) \right] =$$

$$\operatorname{Tr}(\mathbf{H}^{(2;m_2)}^{\top} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}) \operatorname{Tr}(\mathbf{H}^{(1;m_1)}^{\top} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}).$$
(10)

Interpretation. The factorization shows that preserving local geometry among fibers indexed by (i_1,i_2) decomposes into two preservation tasks, one per mode. Each trace term enforces neighborhood smoothness along its mode, while the Kronecker structure captures their joint effect. This regularization encourages nearby slices in the tensor to share similar coefficients in the reduced space, aligning feature selection with the data manifold. The mode-wise decomposition also lowers computational cost and clarifies each mode's contribution to locality preservation.

Similarity Construction. The similarity matrices $\mathbf{B}^{(1)} = [b_{i_1,i_2}^{(1)}] \in \mathbb{R}^{I_1 \times I_1}$ and $\mathbf{B}^{(2)} = [b_{i_1,i_2}^{(2)}] \in \mathbb{R}^{I_2 \times I_2}$ are built via a heat kernel. For example, the similarity between two mode-n slices $\mathbf{X}_{i_1}^{(n)}$ and $\mathbf{X}_{i_2}^{(n)}$, where $n \in \{1,2\}$, is defined as: $b_{i_1,i_2}^{(n)} = \exp\left(-\|\mathbf{X}_{i_1}^{(n)} - \mathbf{X}_{i_2}^{(n)}\|_F^2/\sigma^2\right)$ if $\mathbf{X}_{i_1}^{(n)} \in \mathcal{N}_k(\mathbf{X}_{i_2}^{(n)})$ or vice versa; otherwise $b_{i_1,i_2}^{(n)} = 0$, where $\sigma > 0$ is the kernel width and $\mathcal{N}_k(\cdot)$ denotes the set of k nearest neighbors.

Overall Objective Function. Bringing together the reconstruction fidelity, sparsity control, and manifold preservation, the MSLFS framework can be formulated as

$$\min_{\mathbf{H}^{(n;m_n)}, \mathbf{W}^{(n;m_n)} \geq 0, \forall n \in \{1,2\}} \frac{1}{2} \| \mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2
+ \frac{\alpha}{2} \operatorname{Tr}(\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}) \operatorname{Tr}(\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)})
+ \frac{\beta}{2} \operatorname{Tr}(\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}) \operatorname{Tr}(\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}})
\text{s.t.} \quad \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \otimes \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} = \mathbf{I}_{m_1 m_2}. \quad (11)$$

Details of the optimization procedure, convergence analysis, and computational complexity are provided in Appendices 7.3, 7.4, and 7.5, respectively. In brief, Algorithm 1 outlines the optimization steps for solving the minimization problem (11).

Algorithm 1 MSLFS Algorithm

Input: Data tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$; numbers of selected slices m_1, m_2 ; parameters α, β, γ ; max_iter. Output: Compute ℓ_2 -norm of columns in $\mathbf{W}^{(1;m_1)}$, $\mathbf{W}^{(2;m_2)}$, sort descending. Select top m_1 columns of $\mathbf{W}^{(1;m_1)}$, top m_2 of $\mathbf{W}^{(2;m_2)}$ for mode-1, mode-2 slices. Output $m_1 \times m_2$ features at their intersection.

1: Initialize $\mathbf{W}^{(1;m_1)} \in \mathbb{R}^{m_1 \times I_1}$, $\mathbf{W}^{(2;m_2)} \in \mathbb{R}^{m_2 \times I_2}$, $\mathbf{H}^{(1;m_1)} \in \mathbb{R}^{I_1 \times m_1}$, $\mathbf{H}^{(2;m_2)} \in \mathbb{R}^{I_2 \times m_2}$ randomly; build similarity matrices $\mathbf{B}^{(1)}$, $\mathbf{B}^{(2)}$.

2: for t=0 to max_iter do

3: Update $\mathbf{W}^{(1;m_1)}$ via (17), $\mathbf{H}^{(1;m_1)}$ via (19), $\mathbf{W}^{(2;m_2)}$ via (23), $\mathbf{H}^{(2;m_2)}$ via (25).

5 EXPERIMENTS

In this section, we demonstrate the effectiveness of MSLFS through extensive experiments, comparing it with top-performing feature selection models on real-world benchmark datasets.

Datasets and Compared Methods. To evaluate the effectiveness of MSLFS, we conduct experiments on several benchmark datasets, including **COIL20** (Nene et al., 1996), **ORL** (Cai et al., 2010), **UMIST** (Graham & Allinson, 1998), **Pixraw10P** (Li et al., 2017), **Orlraws10P** (Li et al., 2017),

FashionMNIST (Xiao et al., 2017), BreastMNIST (Yang et al., 2021), and OrganSMNIST (Yang et al., 2021). For comparison, we select ten top-tier models: LS (He et al., 2005b), UDFS (Yang et al., 2011), ILFS (Roffo et al., 2017), GRLTR (Su et al., 2018), CAE (Balın et al., 2019), FSPCA (Tian et al., 2020), CPUFS (Chen et al., 2023), SPCAFS (Li et al., 2023), GRSSLFS (Tiwari et al., 2024), and SPDFS (Dong et al., 2025).

Experimental Settings. To ensure fair evaluation, all methods are tuned under comparable settings. For graph-based approaches, the k-neighborhood is selected from $\{2,5,10,15\}$. We fix $\gamma=10^8$ to enforce orthogonality and set the kernel width $\sigma=10^3$. Regularization parameters are searched over $\{10^{-4},10^{-3},\ldots,10^4\}$, and the number of selected features is varied across $\{50,100,150,200,250,300\}$. Clustering is performed with the true number of clusters, and the maximum iterations of iterative methods are tuned within $\{5,10,30\}$, where 5 or 10 iterations offer a good trade-off between efficiency and convergence. k-means is applied to the selected features and repeated 10 times with random initializations; average results are reported. Performance is assessed by ACC and NMI (Solorio-Fernández et al., 2020), where higher values indicate better results.

Clustering Results. Table 2 presents ACC and NMI results across eight benchmarks against 10 leading baselines. MSLFS consistently achieves top performance, with large improvements on COIL20, ORL, and Orlraws10P, and robust results on challenging datasets such as FashionMNIST and BreastMNIST. These gains come from its slice-based subspace modeling, which leverages cross-mode structure, and graph-regularized selection, which maintains local geometry, producing compact and discriminative features that drive clustering accuracy.

Table 2: Clustering results of the MSLFS vs. 10 cutting-edge models on benchmark datasets.

Model	COIL20		ORL		UMIST		Pixraw10P		Orlraws10P		FashionMNIST		BreastMNIST		OrganSMNIST	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
LS (NeurIPS 2005)	54.34	72.11	48.11	71.44	41.05	59.31	67.37	83.52	69.24	78.76	50.49	51.17	62.22	4.78	33.17	37.69
UDFS (IJCAI 2011)	55.47	71.19	47.95	71.50	36.52	53.03	70.54	79.94	57.64	67.57	52.46	51.22	62.67	5.55	33.52	37.34
ILFS (CVPR 2017)	61.45	73.56	56.68	75.92	45.52	58.74	73.29	83.74	74.52	82.26	63.57	60.31	63.57	7.43	28.86	34.58
GRLTR (JVCIR 2018)	68.78	77.84	54.32	75.00	49.68	63.21	92.44	93.67	82.90	87.51	54.92	51.01	59.19	5.00	33.38	32.16
CAE (ICML 2019)	59.93	72.17	56.25	74.93	54.34	69.22	86.27	91.75	74.45	81.23	67.57	64.26	74.88	9.36	39.81	41.96
FSPCA (NeurIPS 2020)	67.14	79.43	57.07	73.97	52.38	65.54	85.66	92.16	80.41	87.74	63.26	61.68	71.42	8.55	38.15	40.81
CPUFS (TPAMI 2022)	64.72	76.21	57.38	75.39	49.46	63.37	77.27	89.40	76.81	85.36	60.53	58.52	67.87	8.26	37.24	39.57
SPCAFS (TPAMI 2023)	63.15	74.74	52.21	71.76	44.23	58.21	82.16	88.91	73.36	80.44	54.36	51.53	60.46	5.42	34.01	33.26
GRSSLFS (TMLR 2024)	67.47	78.76	53.95	74.58	58.06	68.06	89.30	92.17	79.10	86.04	56.65	62.43	53.85	10.00	32.74	30.94
SPDFS (TPAMI 2025)	67.66	78.96	53.64	73.01	48.37	61.15	78.36	89.13	75.45	82.21	56.76	52.96	61.12	7.66	33.41	34.25
MSLFS (Ours)	73.15	84.67	64.43	79.61	56.79	70.17	93.16	94.28	88.33	91.42	66.42	66.74	76.93	12.85	44.25	44.87
Improvement	+4.37	+5.24	+7.05	+3.69	-	+0.95	+0.72	+0.61	+5.43	+3.68	-	+2.48	+2.05	+2.85	+4.44	+2.91

Ablation Study. The MSLFS objective includes two regularizations: locality preservation (α) to capture local geometry and sparsity (β) to enhance discriminability. An ablation study on six datasets (Table 3) shows that the full model consistently outperforms reduced variants. Removing either term lowers performance, with the sharpest drop when both are omitted, confirming their complementary importance for robust clustering.

Table 3: Ablation study results on six datasets.

Case	COIL20		COIL20 Pixraw10P		ORL		BreastMNIST		UMIST		OrganSMNIST	
Cusc	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
$\alpha, \beta \neq 0$	72.88	83.44	93.66	94.11	64.13	79.45	76.13	12.69	56.33	69.79	44.21	44.57
$\alpha = 0$	66.87	77.22	89.80	91.37	61.12	77.54	63.14	8.44	49.42	63.38	38.54	39.11
$\beta = 0$	68.13	78.97	90.45	92.32	58.98	75.66	68.73	10.89	49.01	61.88	41.11	42.77
$\alpha, \beta = 0$	64.86	74.78	85.20	88.03	56.90	74.05	61.13	7.89	46.56	57.22	36.18	36.66

Convergence Curves. This section analyzes the convergence of MSLFS on four benchmark datasets. Figure 2 shows objective values versus iterations (up to 50). In all cases, the loss drops quickly at first and then stabilizes, demonstrating fast and robust convergence across diverse datasets.



Figure 2: Convergence curves of the MSLFS on the image datasets.

Computational Complexity. Table 4 shows that while many methods incur cubic costs in tensor dimensions, MSLFS reduces complexity to linear dependence on $I_1I_2I_3$ with only minor contributions from slice counts. Its mode-wise design distributes selection across modes and avoids costly global operations, yielding clear efficiency gains over prior approaches.

Table 4: Computational complexity of different models for each iteration. Here t and c denote the dimension of the reduced space and cluster number, respectively.

Model	Computational Complexity
LS	$\mathcal{O}(I_1I_2I_3^2 + I_1I_2\log_2 I_1I_2)$
UDFS	$\mathcal{O}(I_1^3 I_2^3 + I_3^2 c)$
FSPCA	$\mathcal{O}(\max\{I_1I_2m_1m_2t, m_1^3m_2^3\} + I_1I_2m_1m_2t)$
CPUFS	$\mathcal{O}((I_1I_3 + I_2I_3)c^2 + (I_1I_2I_3 + I_3^2)c)$
SPCAFS	$\mathcal{O}(I_1^2 I_2^2 (I_3 + I_1 I_2))$
GRSSLFS	$\mathcal{O}ig(I_1^2I_2^2I_3^2ig)$
GRLTR	$\mathcal{O}(I_1I_2I_3\log_2I_3 + I_1I_2^2I_3 + I_3^3)$
SPDFS	$\mathcal{O}\left(\max\{I_1I_2I_3t, I_1^2I_2^2t\} + I_1I_2I_3tc^2 + \max\{I_1I_2I_3t, I_1I_2\log_2I_1I_2, m_1^3m_2^3\}\right)$
MSLFS	$\mathcal{O}(I_1I_2I_3(\max\{m_1,I_2\}+\max\{m_2,I_1\}))$

Data Visualization using t-SNE. Figure 3 presents t-SNE visualizations on UMIST. The raw data shows scattered and overlapping clusters, while MSLFS with varying feature counts produces progressively clearer, more compact, and better-separated groups. This demonstrates the ability of the MSLFS to extract discriminative features that enhance clustering quality.

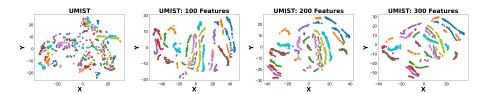


Figure 3: t-SNE plots of UMIST before and after feature reduction with MSLFS.

Selected Features Visualization. Figure 4 depicts feature selection on ORL and Pixraw10P with 100, 200, and 300 features. Fewer features capture broad structure, while more reveal finer details. Across datasets, the model consistently highlights informative regions, expressing its efficacy for image-based feature selection.

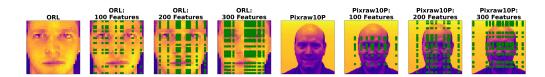


Figure 4: Image visualizations on ORL and Pixraw10P with 100, 200, and 300 selected features.

6 CONCLUSION

The proposed MSLFS introduces a novel approach to tensor-based feature selection by distributing the selection process across modes rather than treating the feature space as a rigid whole. Its key innovation, the multi-linear subspace distance, provides a principled criterion for preserving global structure while enabling efficient and interpretable feature selection. Complemented by joint sparsity and higher-order graph regularization, MSLFS captures both cross-mode dependencies and local manifold geometry, setting it apart from existing tensor-based methods. This framework opens new directions for multi-way learning, with future work aimed at extending MSLFS to broader tasks such as its integration with deep tensor architectures for large-scale representation learning. Comprehensive theoretical discussions and supplementary experiments can be found in Appendices 7 and 8*.

^{*}This paper has used large language models solely for improving the clarity and polish of the writing.

REFERENCES

- Muhammed Fatih Balin, Abubakar Abid, and James Zou. Concrete Autoencoders: Differentiable Feature Selection and Reconstruction. In *International Conference on Machine Learning*, pp. 444–453. PMLR, 2019.
- Qi Bi, Jingjun Yi, Hao Zheng, Wei Ji, Yawen Huang, Yuexiang Li, and Yefeng Zheng. Learning Generalized Medical Image Representation by Decoupled Feature Queries. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–18, 2025.
- Deng Cai, Chiyuan Zhang, and Xiaofei He. Unsupervised Feature Selection for Multi-Cluster Data. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 333–342, 2010.
- Bilian Chen, Jiewen Guan, and Zhening Li. Unsupervised Feature Selection via Graph Regularized Non-Negative CP Decomposition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):2582–2594, 2023.
- Ammar Chouchane, Mohcene Bessaoudi, Hamza Kheddar, Abdelmalik Ouamane, Tiago Vieira, and Mahmoud Hassaballah. Multilinear Subspace Learning for Person Re-Identification Based Fusion of High Order Tensor Features. *Engineering Applications of Artificial Intelligence*, 128: 107521, 2024.
- Xia Dong, Feiping Nie, Lai Tian, Rong Wang, and Xuelong Li. Unsupervised Discriminative Feature Selection with $\ell_{2,0}$ -Norm Constrained Sparse Projection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–15, 2025.
- Daniel B Graham and Nigel M Allinson. Characterising Virtual Eigensignatures for General Purpose Face Recognition. In *Face Recognition: From Theory to Applications*, pp. 446–456. Springer, 1998.
- Xiaofei He, Deng Cai, and Partha Niyogi. Laplacian Score for Feature Selection. *Advances in Neural Information Processing Systems*, 18, 2005a.
- Xiaofei He, Deng Cai, and Partha Niyogi. Laplacian Score for Feature Selection. *Advances in Neural Information Processing Systems*, 18, 2005b.
- Hamed Jelodar, Yongli Wang, Chi Yuan, Xia Feng, Xiahui Jiang, Yanchao Li, and Liang Zhao. Latent Dirichlet Allocation (LDA) and Topic Modeling: Models, Applications, a Survey. *Multimedia tools and applications*, 78(11):15169–15211, 2019.
- Daniel D. Lee and H. Sebastian Seung. Learning the Parts of Objects by Non-Negative Matrix Factorization. *Nature*, 401:788–791, 10 1999.
- Hua Li, Wenya Luo, Zhidong Bai, Huanchao Zhou, and Zhangni Pu. Spectrally-Corrected and Regularized LDA for Spiked Model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(3):1991–1999, 2025.
- Jundong Li, Kewei Cheng, Suhang Wang, Fred Morstatter, Robert P Trevino, Jiliang Tang, and Huan Liu. Feature Selection: A Data Perspective. *ACM Computing Surveys (CSUR)*, 50(6):1–45, 2017.
- Zhengxin Li, Feiping Nie, Jintang Bian, Danyang Wu, and Xuelong Li. Sparse PCA via $\ell_{2,p}$ -Norm Regularization for Unsupervised Feature Selection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):5322–5328, 2023.
- Zhi Liu, Bo Tang, Xiaofu He, Qingchen Qiu, and Hongjun Wang. Sparse Tensor-Based Dimensionality Reduction for Hyperspectral Spectral–Spatial Discriminant Feature Extraction. *IEEE Geoscience and Remote Sensing Letters*, 14(10):1775–1779, 2017.
- Canyi Lu, Jiashi Feng, Yudong Chen, Wei Liu, Zhouchen Lin, and Shuicheng Yan. Tensor Robust Principal Component Analysis with a New Tensor Nuclear Norm. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):925–938, 2020.

- Haiping Lu, Konstantinos N. Plataniotis, and Anastasios N. Venetsanopoulos. A Survey of Multi-linear Subspace Learning for Tensor Data. *Pattern Recognition*, 44(7):1540–1551, 2011.
 - Sameer A Nene, Shree K Nayar, Hiroshi Murase, et al. Columbia Object Image Library (COIL-100). Technical report, Technical Report CUCS-006-96, 1996.
 - Giorgio Roffo, Simone Melzi, Umberto Castellani, and Alessandro Vinciarelli. Infinite Latent Feature Selection: A Probabilistic Latent Graph-Based Ranking Approach. In *Proceedings of the IEEE international conference on computer vision*, pp. 1398–1406, 2017.
 - Saúl Solorio-Fernández, J Ariel Carrasco-Ochoa, and José Fco Martínez-Trinidad. A Review of Unsupervised Feature Selection Methods. *Artificial Intelligence Review*, 53(2):907–948, 2020.
 - Kun Song, Hao Li, Gong Cheng, Junwei Han, Feiping Nie, Bin Gu, and Fakhri Karray. Learning Compact Discriminant Representation via Low-rank Bilinear Pooling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–18, 2025.
 - Yuting Su, Xu Bai, Wu Li, Peiguang Jing, Jing Zhang, and Jing Liu. Graph Regularized Low-Rank Tensor Representation for Feature Selection. *Journal of Visual Communication and Image Representation*, 56:234–244, 2018.
 - Lai Tian, Feiping Nie, Rong Wang, and Xuelong Li. Learning Feature Sparse Principal Subspace. *Advances in Neural Information Processing Systems*, 33:14997–15008, 2020.
 - Prayag Tiwari, Farid Saberi Movahed, Saeed Karami, Farshad Saberi-Movahed, Jens Lehmann, and Sahar Vahdati. A Self-Representation Learning Method for Unsupervised Feature Selection using Feature Space Basis. *Transactions on Machine Learning Research*, 2024.
 - Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: A Novel Image Dataset for Benchmarking Machine Learning Algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
 - Jiancheng Yang, Rui Shi, and Bingbing Ni. MedMNIST Classification Decathlon: A Lightweight Auto-ML Benchmark for Medical Image Analysis. In 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI), pp. 191–195. IEEE, 2021.
 - Yi Yang, Heng Tao Shen, Zhigang Ma, Zi Huang, and Xiaofang Zhou. $\ell_{2,1}$ -Norm Regularized Discriminative Feature Selection for Unsupervised Learning. In *IJCAI International Joint Conference on Artificial Intelligence*, 2011.
 - Ron Zass and Amnon Shashua. Non-Negative Sparse PCA. Advances in Neural Information Processing Systems, 19, 2006.

7 ADDITIONAL THEORETICAL RESULTS

7.1 NOTATIONS AND PRELIMINARIES

Notations. Throughout this paper, vectors are represented by bold lowercase letters (e.g., \mathbf{v}), matrices by bold uppercase letters (e.g., \mathbf{A}), and tensors by bold calligraphy letters (e.g., \mathcal{X}). The identity matrix of size m is denoted by \mathbf{I}_m , and $\mathbf{e}_j^{(m)}$ denotes the jth column of \mathbf{I}_m . For a matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{m \times n}$, the i-th row and the j-th column are denoted by $\mathbf{A}_{i,:}$ and $\mathbf{A}_{:,j}$, respectively. The Frobenius norm of \mathbf{A} is defined as $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$, while the $\ell_{2,1}$ -norm is given by $\|\mathbf{A}\|_{2,1} = \sum_{i=1}^m \|\mathbf{A}_{i,:}\|_2$. For a square matrix \mathbf{A} , $\mathrm{Tr}(\mathbf{A})$ denotes its trace. The dot product between two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ is defined as $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^n u_i v_i$, and the Frobenius inner product between two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is defined as $\langle \mathbf{A}, \mathbf{B} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n a_{ij} b_{ij}$. The Hadamard product of $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is expressed as $\mathbf{A} \odot \mathbf{B} = (a_{ij}b_{ij})_{i=1,j=1}^{m,n}$. For $\mathbf{A} \in \mathbb{K}^{I \times J}$ and $\mathbf{B} \in \mathbb{K}^{M \times N}$, the Kronecker product is $\mathbf{A} \otimes \mathbf{B} \in \mathbb{K}^{(IM) \times (JN)}$; defining $\overline{im} = (i-1)M + m$ and $\overline{jn} = (j-1)N + n$, its entries are given by $(\mathbf{A} \otimes \mathbf{B})_{\overline{im},\overline{jn}} = a_{ij}b_{mn}$ for $1 \leq i \leq I, 1 \leq j \leq J, 1 \leq m \leq M$, and $1 \leq n \leq N$. A third-order tensor is denoted by $\mathcal{X} = (x_{i_1,i_2,i_3})_{i_1=1,\dots,I_1;\ i_2=1,\dots,I_2;\ i_3=1,\dots,I_3} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, where each entry $x_{i_1,i_2,i_3} \in \mathbb{R}$; the j-th frontal slice $(j=1,\dots,I_3)$, denoted by $\mathbf{X}_j^{(3)}$, is obtained by fixing the third index and belongs to $\mathbb{R}^{I_1 \times I_2}$. Furthermore, the mode-3 unfolding of \mathcal{X} , denoted by $\mathbf{X}_{(3)}$, rearranges the entries of \mathcal{X} into a matrix of size $I_3 \times I_1 I_2$ by mapping the mode-3 fibers to the columns of $\mathbf{X}_{(3)}$.

Preliminaries. Let $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$ be an N-mode tensor, $\mathbf{A} \in \mathbb{R}^{J \times I_n}$ a matrix, and $\mathbf{v} \in \mathbb{R}^{I_n}$ a vector $(n=1,\ldots,N)$. The n-mode tensor-matrix product $\mathcal{X} \times_n \mathbf{A} \in \mathbb{R}^{I_1 \times \cdots \times J_N \cdots \times I_N}$ and the n-mode tensor-vector product $\mathcal{X} \bar{\times}_n \mathbf{v} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times I_{n+1} \times \cdots \times I_N}$ are defined elementwise as $(\mathcal{X} \times_n \mathbf{A})_{i_1 \cdots j_1 \cdots i_N} = \sum_{i_n=1}^{I_n} x_{i_1 \cdots i_N} a_{ji_n}, \quad (\mathcal{X} \bar{\times}_n \mathbf{v})_{i_1 \cdots i_{n-1} i_{n+1} \cdots i_N} = \sum_{i_n=1}^{I_n} x_{i_1 \cdots i_N} v_{i_n}.$ For $\mathbf{y} \in \mathbb{R}^J$, we have $\mathcal{X} \times_n \mathbf{A} \bar{\times}_n \mathbf{y} = \mathcal{X} \bar{\times}_n (\mathbf{y}^\top \mathbf{A})$, and the j-th mode-n slice of $\mathcal{X} \times_n \mathbf{A}$ is $\mathcal{X} \bar{\times}_n \mathbf{A}_{j,:}$, $j=1,\ldots,J$. In particular, if $\mathbf{y}=\mathbf{e}_j^{(n)\top}$ is the j-th column of the identity \mathbf{I}_{I_n} , then $\mathcal{X} \times_n \mathbf{y}$ extracts the j-th mode-n slice: $\mathcal{X} \times_n \mathbf{e}_j^{(n)\top} = \mathbf{X}_j^{(n)}$.

The mode-n unfolding of \mathcal{X} , denoted by $\mathbf{X}_{(n)} \in \mathbb{R}^{I_n \times (I_1 \cdots I_{n-1}I_{n+1} \cdots I_N)}$, rearranges \mathcal{X} into a matrix by aligning all mode-n fibers as its columns. Tensor-matrix products admit the following unfolding formulations: $(\mathcal{X} \times_n \mathbf{A})_{(n)} = \mathbf{A} \, \mathbf{X}_{(n)}, \, (\mathcal{X} \times_m \mathbf{A})_{(n)} = \mathbf{X}_{(n)} \big(\mathbf{I}_{I_{m+1} \cdots I_N} \otimes \mathbf{A} \otimes \mathbf{I}_{I_1 \cdots I_{m-1}} \big)^\top, m \neq n$. More generally, for a sequence of tensor-matrix products, we have $(\mathcal{X} \times_1 \mathbf{A}_1 \times_2 \mathbf{A}_2 \cdots \times_N \mathbf{A}_N)_{(n)} = \mathbf{A}_n \, \mathbf{X}_{(n)} \Big(\mathbf{A}_N \otimes \cdots \otimes \mathbf{A}_{n+1} \otimes \mathbf{A}_{n-1} \otimes \cdots \otimes \mathbf{A}_1 \Big)^\top$.

7.2 Proof of Theorem 3.1

Proof. Step 1 (Core matrix equals intersection fibers). Define the subtensor

$$\mathcal{Y} := \mathcal{X} \times_1 \mathbf{W}^{(1;R_1)} \times_2 \mathbf{W}^{(2;R_2)} \in \mathbb{R}^{R_1 \times R_2 \times I_3}.$$

Since $\mathbf{W}^{(1;R_1)}$ and $\mathbf{W}^{(2;R_2)}$ are indicator selectors for the index sets T_1 and T_2 , the $(r_1,r_2,:)$ -entry of $\mathcal Y$ is precisely the intersection fiber $\mathbf{f}_{i^{(1)}_{r_1},i^{(2)}_{r_2}}=\mathcal X_{i^{(1)}_{r_1},i^{(2)}_{r_2}}$. Hence, $\mathcal Y$ stacks exactly the R_1R_2 intersection fibers. Unfolding along mode-3 and using the standard product—unfolding identity yields

$$\mathbf{F}_{\mathrm{core}} := \mathbf{Y}_{(3)} = \mathbf{X}_{(3)} (\mathbf{W}^{(2;R_2)} \otimes \mathbf{W}^{(1;R_1)})^\top,$$

so \mathbf{F}_{core} is exactly the matrix whose columns are the R_1R_2 intersection fibers.

Step 2 (Coefficient matrices along modes 1 and 2). Because $\{\mathbf{X}_{i_{r_1}^{(1)}}^{(1)}\}_{r_1=1}^{R_1}$ is a basis of $\mathcal{S}(\mathbf{X}^{(1)})$,

for every $i_1 \in \{1,2,\cdots,I_1\}$, there exist coefficients $\{h_{i_1,r_1}^{(1;R_1)}\}_{r_1=1}^{R_1}$ such that

$$\mathbf{X}_{i_1}^{(1)} = \sum_{r_1=1}^{R_1} h_{i_1,r_1}^{(1;R_1)} \, \mathbf{X}_{i_{r_1}^{(1)}}^{(1)}.$$

Collect these into $\mathbf{H}^{(1;R_1)} \in \mathbb{R}^{I_1 \times R_1}$. Similarly, since $\{\mathbf{X}_{i_{r_2}}^{(2)}\}_{r_2=1}^{R_2}$ is a basis of $\mathcal{S}(\mathbf{X}^{(2)})$, there exists $\mathbf{H}^{(2;R_2)} \in \mathbb{R}^{I_2 \times R_2}$ with

$$\mathbf{X}_{i_2}^{(2)} = \sum_{r_2=1}^{R_2} h_{i_2, r_2}^{(2; R_2)} \, \mathbf{X}_{i_{r_2}^{(2)}}^{(2)}.$$

Step 3 (Fiber-level decomposition). Fix (i_1, i_2) . Expanding along mode 1 and then mode 2 gives

$$\mathbf{f}_{i_1,i_2} = \mathcal{X}_{i_1,i_2,:} = \sum_{r_1=1}^{R_1} h_{i_1,r_1}^{(1;R_1)} \, \mathcal{X}_{i_{r_1}^{(1)},i_2,:} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} h_{i_1,r_1}^{(1;R_1)} h_{i_2,r_2}^{(2;R_2)} \, \mathcal{X}_{i_{r_1}^{(1)},i_{r_2}^{(2)},:}$$

i.e.,

$$\mathbf{f}_{i_1,i_2} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} h_{i_1,r_1}^{(1;R_1)} h_{i_2,r_2}^{(2;R_2)} \, \mathbf{f}_{i_{r_1}^{(1)},i_{r_2}^{(2)}}.$$

Step 4 (Tensor-level identity). Stacking the identities in Step 3 over all (i_1, i_2) shows that \mathcal{X} is obtained by first selecting the basis slices and then recombining them with the coefficients:

$$\mathcal{X} = \mathcal{X} \times_1 \mathbf{H}^{(1;R_1)} \mathbf{W}^{(1;R_1)} \times_2 \mathbf{H}^{(2;R_2)} \mathbf{W}^{(2;R_2)}.$$

Unfolding this equality along mode 3 and using the same product-unfolding identity as in Step 1 gives

$$\begin{split} \mathbf{X}_{(3)} &= \left(\mathcal{X} \times_1 \mathbf{H}^{(1;R_1)} \mathbf{W}^{(1;R_1)} \times_2 \mathbf{H}^{(2;R_2)} \mathbf{W}^{(2;R_2)} \right)_{(3)} \\ &= \mathbf{X}_{(3)} \left(\mathbf{W}^{(2;R_2)} \otimes \mathbf{W}^{(1;R_1)} \right)^{\top} \left(\mathbf{H}^{(2;R_2)} \otimes \mathbf{H}^{(1;R_1)} \right)^{\top}. \end{split}$$

Substituting the expression for \mathbf{F}_{core} from Step 1 yields

$$\mathbf{X}_{(3)} = \mathbf{F}_{\text{core}}(\mathbf{H}^{(2;R_2)} \otimes \mathbf{H}^{(1;R_1)})^{\top}.$$

This completes the proof.

7.3 OPTIMIZATION

We now detail the optimization procedure of the proposed MSLFS method given in Problem (12), describing the iterative steps for solving its objective function and updating the associated optimization variables.

$$\min_{\mathbf{H}^{(n;m_n)}, \mathbf{W}^{(n;m_n)} \geq 0, \forall n \in \{1,2\}} \frac{1}{2} \| \mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2
+ \frac{\alpha}{2} \operatorname{Tr}(\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}) \operatorname{Tr}(\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)})
+ \frac{\beta}{2} \operatorname{Tr}(\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}) \operatorname{Tr}(\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}})
\text{s.t.} \quad \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \otimes \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} = \mathbf{I}_{m_1 m_2}. \quad (12)$$

To derive the multiplicative updating rules for $\mathbf{W}^{(1;m_1)}$ and $\mathbf{H}^{(1;m_1)}$, one must calculate the derivatives of the objective function with respect to these variables and set them equal to zero. To this end, the first term of the objective function can be unfolded as follows:

$$\min_{\mathbf{H}^{(n;m_n)}, \mathbf{W}^{(n;m_n)} \geq 0, \forall n \in \{1,2\}} \frac{1}{2} \| \mathbf{X}_{(1)} - \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \mathbf{X}_{(1)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)})^{\top} \|_F^2
+ \frac{\alpha}{2} \mathrm{Tr} [\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \mathrm{Tr} [\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}]
+ \frac{\beta}{2} \mathrm{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}] \mathrm{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}}]
\text{s.t.} \quad \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \otimes \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} = \mathbf{I}_{m_1 m_2}. \quad (13)$$

Simplifying the objective function leads us to:

$$\frac{1}{2} \text{Tr}[\mathbf{X}_{(1)}^{\top} \mathbf{X}_{(1)} - 2(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(1;m_{1})^{\top}} \mathbf{H}^{(1;m_{1})^{\top}} \mathbf{X}_{(1)} \\
+ (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(1;m_{1})^{\top}} \mathbf{H}^{(1;m_{1})^{\top}} \mathbf{H}^{(1;m_{1})} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})})^{\top}] \\
+ \frac{\alpha}{2} \text{Tr}[\mathbf{H}^{(2;m_{2})^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_{2})}] \text{Tr}[\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_{1})}] \\
+ \frac{\beta}{2} \text{Tr}[\mathbf{W}^{(2;m_{2})} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_{2})^{\top}}] \text{Tr}[\mathbf{W}^{(1;m_{1})} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_{1})^{\top}}] \\
+ \frac{\gamma}{4} \text{Tr}[\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}} \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}} \otimes \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}} \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}} \\
- 2 \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}} \otimes \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}} + \mathbf{I}_{m_{1}m_{2}}]. \tag{14}$$

Now the the derivatives of the objective function w.r.t. $\mathbf{W}^{(1;m_1)}$ and $\mathbf{H}^{(1;m_1)}$ can be calculated as follows:

$$\frac{\partial \mathcal{F}}{\partial \mathbf{W}^{(1;m_1)}} = -\mathbf{H}^{(1;m_1)^{\top}} \mathbf{X}_{(1)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)}) \mathbf{X}_{(1)}^{\top}
+ \mathbf{H}^{(1;m_1)^{\top}} \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \mathbf{X}_{(1)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)})^{\top} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)}) \mathbf{X}_{(1)}^{\top}
+ \beta \text{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}] \mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)}
+ \gamma \text{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}}] \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)}
- \gamma \text{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}}] \mathbf{W}^{(1;m_1)}.$$
(15)

$$\frac{\partial \mathcal{F}}{\partial \mathbf{H}^{(1;m_1)}} = -\mathbf{X}_{(1)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(1;m_1)^{\top}}
+ \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \mathbf{X}_{(1)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)})^{\top} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(1;m_1)^{\top}}
+ \alpha \text{Tr} [\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}.$$
(16)

According to KKT conditions (Lee & Seung, 1999), we have the following updating rules:

$$\mathbf{W}^{(1;m_{1})} = (\mathbf{W}^{(1;m_{1})} \odot \mathbf{H}^{(1;m_{1})^{\top}} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top}$$

$$+ \gamma \text{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}] \mathbf{W}^{(1;m_{1})})$$

$$\otimes (\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top}$$

$$+ \beta \text{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_{2})^{\top}}] \mathbf{W}^{(1;m_{1})} \mathbf{U}^{(1)}$$

$$+ \gamma \text{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}} \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}] \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}} \mathbf{W}^{(1;m_{1})}$$

$$+ \gamma \text{Tr} [\mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(1;m_{1})^{\top}}$$

$$+ \alpha \text{Tr} [\mathbf{H}^{(2;m_{2})^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_{2})}] \mathbf{A}^{(1)} \mathbf{H}_{1})$$

$$\otimes (\mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(2;m_{2})}$$

$$+ \alpha \text{Tr} [\mathbf{H}^{(2;m_{2})^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_{2})}] \mathbf{B}^{(1)} \mathbf{H}^{(1;m_{1})}).$$
(19)

To derive the update rules for $\mathbf{W}^{(2;m_2)}$ and $\mathbf{H}^{(2;m_2)}$, the first term of (12) must be reformulated using the mode-2 unfolding of the tensor. Then, the derivatives of (12) with respect to these variables are computed.

$$\min_{\mathbf{H}^{(n;m_n)}, \mathbf{W}^{(n;m_n)} \geq 0, \forall n \in \{1,2\}} \frac{1}{2} \| \mathbf{X}_{(2)} - \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \mathbf{X}_{(2)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)})^{\top} \|_F^2
+ \frac{\alpha}{2} \text{Tr} [\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \text{Tr} [\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}]
+ \frac{\beta}{2} \text{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}] \text{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}}]
s.t. \quad \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \otimes \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} = \mathbf{I}_{m_1 m_2}. \quad (20)$$

$$\frac{\partial \mathcal{F}}{\partial \mathbf{W}^{(2;m_2)}} = -\mathbf{H}^{(2;m_2)^{\top}} \mathbf{X}_{(2)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)}) \mathbf{X}_{(2)}^{\top}
+ \mathbf{H}^{(2;m_2)^{\top}} \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \mathbf{X}_{(2)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)})^{\top} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)}) \mathbf{X}_{(2)}^{\top}
+ \beta \text{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}}] \mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)}
+ \gamma \text{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}}] \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \mathbf{W}^{(2;m_2)}
- \gamma \text{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}}] \mathbf{W}^{(2;m_2)}.$$
(21)

$$\frac{\partial \mathcal{F}}{\partial \mathbf{H}^{(2;m_2)}} = -\mathbf{X}_{(2)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_2)^{\top}}
+ \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \mathbf{X}_{(2)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)})^{\top} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_2)^{\top}}
+ \alpha \text{Tr} [\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}] \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}.$$
(22)

According to KKT conditions (Lee & Seung, 1999), we have the following updating rules:

$$\mathbf{W}^{(2;m_{2})} = \mathbf{W}^{(2;m_{2})} \odot (\mathbf{H}^{(2;m_{2})^{\top}} \mathbf{X}_{(2)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top}$$

$$+ \gamma \text{Tr} [\mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}}] \mathbf{W}^{(2;m_{2})})$$

$$\otimes (\mathbf{H}^{(2;m_{2})^{\top}} \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})} \mathbf{X}_{(2)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top}$$

$$+ \beta \text{Tr} [\mathbf{W}^{(1;m_{1})} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_{1})^{\top}}] \mathbf{W}^{(2;m_{2})} \mathbf{U}^{(2)}$$

$$+ \gamma \text{Tr} [\mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}} \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}}] \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}),$$
(24)

$$\mathbf{H}^{(2;m_{2})} = \mathbf{H}^{(2;m_{2})} \odot (\mathbf{X}_{(2)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_{2})^{\top}}$$

$$+ \alpha \text{Tr}[\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_{1})}] \mathbf{A}^{(2)} \mathbf{H}^{(2;m_{2})})$$

$$\otimes (\mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})} \mathbf{X}_{(2)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(i;m_{1})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_{2})^{\top}}$$

$$+ \alpha \text{Tr}[\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_{1})}] \mathbf{B}^{(2)} \mathbf{H}^{(2;m_{2})}).$$
(25)

7.4 Convergence Analysis

This section iinvestigates the convergence analysis of MSLFS to explore the decreasing behavior of the objective function (12). It is first assumed that each matrix $\mathbf{W}^{(n;m_n)}$, $\mathbf{H}^{(n;m_n)}$, for $n \in \{1,2\}$ is individually updated while the others remain unchanged. Based on this assumption, the decreasing behavior of the objective function is analyzed for each variable. For this purpose, several important definitions and findings from (Lee & Seung, 1999) are examined.

Definition 7.1 ((Lee & Seung, 1999)). The function $G(u, u^{(t)})$ is deemed an auxiliary function for f(u) if it fulfills the subsequent criteria:

$$g(u, u^{(t)}) \ge f(u), \quad g(u, u) = f(u),$$
 (26)

for every $u \in \mathbb{R}$.

Lemma 1 ((Lee & Seung, 1999)). Suppose $g(u, u^{(t)})$ is an auxiliary function associated with f(u). Then, the sequence $\{f(u^{(t)})\}_{t=1}^{\infty}$ is non-increasing when u is updated according to

$$u^{(t+1)} = \arg\min_{u \in \mathbb{R}} g(u, u^{(t)}).$$

In Proposition 7.2, an auxiliary function is created to ensure that the original objective function diminishes monotonically in line with the update rule for $\mathbf{W}^{(1;m_1)}$ specified in (17).

Proposition 7.2. Given that the matrices $\mathbf{H}^{(1;m_1)}$, $\mathbf{W}^{(2;m_2)}$, and $\mathbf{H}^{(2;m_2)}$ are fixed, the update rule (17) for $\mathbf{W}^{(1;m_1)}$ ensures that the objective function of the minimization problem (12) does not increase.

Proof. Assume that the matrices $\mathbf{H}^{(1;m_1)}$, $\mathbf{W}^{(2;m_2)}$, and $\mathbf{H}^{(2;m_2)}$ are fixed. Consider the objective function in the optimization problem (12) with respect to $\mathbf{W}^{(1;m_1)}$:

$$f(\mathbf{W}^{(1;m_1)}) = \frac{1}{2} \| \mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2$$

$$+ \frac{\beta}{2} \operatorname{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}] \operatorname{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}}]$$

$$+ \frac{\gamma}{4} \| \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \otimes \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} - \mathbf{I}_{m_1 m_2} \|_F^2.$$

To show that $f(\mathbf{W}^{(1;m_1)^{(t+1)}}) \leq f(\mathbf{W}^{(1;m_1)^{(t)}})$, define $g(w_1, f(w_{j_1, i_1}^{(1;m_1)^{(t)}}))$ as follows:

$$g(w_{1}, w_{j_{1}, i_{1}}^{(1;m_{1})(t)}) = \mathcal{B}(w_{j_{1}, i_{1}}^{(1;m_{1})(t)}) + \dot{\mathcal{B}}(w_{j_{1}, i_{1}}^{(1;m_{1})(t)})(w_{1} - w_{j_{1}, i_{1}}^{(1;m_{1})(t)})$$

$$+ \left(\mathbf{H}^{(1;m_{1})^{\top}}\mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})(t)}\mathbf{X}_{(1)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})^{\top}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})\mathbf{X}_{(1)}^{\top}\right)$$

$$+ \beta \text{Tr}[\mathbf{W}^{(2;m_{2})}\mathbf{U}^{(2)}\mathbf{W}^{(2;m_{2})^{\top}}]\mathbf{W}^{(1;m_{1})(t)}\mathbf{U}^{(1)}$$

$$+ \gamma \text{Tr}[\mathbf{W}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})^{\top}}\mathbf{W}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})^{\top}}]\mathbf{W}^{(1;m_{1})(t)}(\mathbf{W}^{(1;m_{1})(t)})^{\top}\mathbf{W}^{(1;m_{1})(t)})$$

$$\times \frac{(w_{1} - w_{j_{1}, i_{1}}^{(1;m_{1})(t)})^{2}}{2 w_{j_{1}, i_{1}}^{(1;m_{1})(t)}},$$

for $j_1=1,2,\cdots,m_1$ and $i_1=1,2,\cdots,I_1$. Moreover, assume that $\mathcal{B}(w_1)$ indicates the part of f(w) relevant to $\mathbf{W}_{j_1,i_1}^{(1;m_1)}$, and

$$\dot{\mathcal{B}}(w_{1}) := \left(\frac{\partial f}{\partial \mathbf{W}^{(1;m_{1})}}\right)_{j_{1},i_{1}} = \left(-\mathbf{H}^{(1;m_{1})^{\top}}\mathbf{X}_{(1)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})\mathbf{X}_{(1)}^{\top} + \mathbf{H}^{(1;m_{1})^{\top}}\mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})}\mathbf{X}_{(1)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})^{\top}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})\mathbf{X}_{(1)}^{\top} + \beta \operatorname{Tr}[\mathbf{W}^{(2;m_{2})}\mathbf{U}^{(2)}\mathbf{W}^{(2;m_{2})^{\top}}]\mathbf{W}^{(1;m_{1})}\mathbf{U}^{(1)} + \gamma \operatorname{Tr}[\mathbf{W}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})^{\top}}\mathbf{W}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})^{\top}}]\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})} - \gamma \operatorname{Tr}[\mathbf{W}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})^{\top}}]\mathbf{W}^{(1;m_{1})}\right)_{j_{1},j_{1}}.$$

It can be seen that $g(w_1, w_{j_1, i_1}^{(1;m_1)^{(t)}})$ is an auxiliary function of $\mathcal{B}(w_1)$. For this purpose, consider the Taylor expansion of $\mathcal{B}(w_1)$ around $w_{j_1, i_1}^{(1;m_1)^{(t)}}$:

$$\mathcal{B}(w_1) = \mathcal{B}(w_{j_1,i_1}^{(1;m_1)^{(t)}}) + \dot{\mathcal{B}}(w_{j_1,i_1}^{(1;m_1)^{(t)}})(w_1 - w_{j_1,i_1}^{(1;m_1)^{(t)}}) + \frac{1}{2}\ddot{\mathcal{B}}(w_{j_1,i_1}^{(1;m_1)^{(t)}})(w_1 - w_{j_1,i_1}^{(1;m_1)^{(t)}})^2,$$

where

$$\ddot{\mathcal{B}}(w_{1}) := \left(\frac{\partial^{2} F}{\partial \mathbf{W}^{(1;m_{1})^{2}}}\right)_{j_{1},i_{1}} = \left(\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{H}^{(1;m_{1})}\right)_{j_{1},j_{1}} \\
\times \left(\mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top}\right)_{i_{1},i_{1}} \\
+ \beta \operatorname{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_{2})^{\top}}] u_{j_{1},i_{1}}^{(1)} \\
+ \gamma \operatorname{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}} \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}] \left(\left(\mathbf{W}^{(1;m_{1})^{\top}} \mathbf{W}^{(1;m_{1})}\right)_{i_{1},i_{1}} + w_{j_{1},i_{1}}^{(1,m_{1})^{2}} \right. \\
+ \left. \left(\mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}}\right)_{j_{1},j_{1}}\right) - \gamma \operatorname{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}].$$

It is easy to validate that $g(w_1, w_1) = \mathcal{B}(w_1)$. Moreover, in light of the following inequalities,

$$\begin{split} & \left(\mathbf{H}^{(1;m_{1})^{\top}}\mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})}\mathbf{X}_{(1)}(\mathbf{I}_{I_{3}}\otimes\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})^{\top}(\mathbf{I}_{I_{3}}\otimes\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})\mathbf{X}_{(1)}^{\top}\right)_{j_{1},i_{1}} \\ & = \sum_{r=1}^{m_{1}}\sum_{s=1}^{I_{1}}\left(\mathbf{H}^{(1;m_{1})^{\top}}\mathbf{H}^{(1;m_{1})}\right)_{j_{1},r}w_{r,s}^{(1;m_{1})}\times\left(\mathbf{X}_{(1)}(\mathbf{I}_{I_{3}}\otimes\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})^{\top}(\mathbf{I}_{I_{3}}\otimes\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})\mathbf{X}_{(1)}^{\top}\right)_{s,i_{1}} \\ & \geq \left(\mathbf{H}^{(1;m_{1})^{\top}}\mathbf{H}^{(1;m_{1})}\right)_{j_{1},j_{1}}\left(\mathbf{X}_{(1)}(\mathbf{I}_{I_{3}}\otimes\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})^{\top}(\mathbf{I}_{I_{3}}\otimes\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})})\mathbf{X}_{(1)}^{\top}\right)_{i_{1},i_{1}}, \\ & \left(\mathbf{W}^{(1;m_{1})}\mathbf{U}^{(1)}\right)_{j_{1},i_{1}} = \sum_{s=1}^{I_{1}}w_{j_{1},s}^{(1,m_{1})}u_{s,i_{1}}^{(1)} \geq u_{j_{1},i_{1}}^{(1)}, \end{split}$$

and

$$\left(\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}\mathbf{W}^{(1;m_{1})}\right)_{j_{1},i_{1}} = \sum_{s=1}^{I_{1}} \sum_{r=1}^{m_{1}} w_{j_{1},s}^{(1,m_{1})} w_{r,s}^{(1,m_{1})} w_{r,i_{1}}^{(1,m_{1})} \ge \left(\mathbf{W}^{(1;m_{1})^{\top}}\mathbf{W}^{(1;m_{1})}\right)_{i_{1},i_{1}} + w_{j_{1},i_{1}}^{(1,m_{1})^{2}} + \left(\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}\right)_{j_{1},j_{1}},$$

it can be observed that $g(w_1, w_{j_1, i_1}^{(1;m_1)^{(t)}}) \geq \mathcal{B}(w_1)$, for each $w_1 \in \mathbb{R}$. Consequently, since the requirements of Definition 26 are met $g(w_1, w_{j_1, i_1}^{(1;m_1)})$ serves as an auxiliary function for $\mathcal{B}(w_1)$. Then, by minimizing $g(w_1, w_{j_1, i_1}^{(1;m_1)^{(t)}})$ with respect to w_1 , the updating rule of $\mathbf{W}^{(1;m_1)}$ can be obtained in the form

$$\mathbf{W}^{(1;m_{1})} = \left(\mathbf{W}^{(1;m_{1})} \odot \mathbf{H}^{(1;m_{1})^{\top}} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top} + \gamma \operatorname{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}] \mathbf{W}^{(1;m_{1})} \right)$$

$$\oslash \left(\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \mathbf{X}_{(1)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})}) \mathbf{X}_{(1)}^{\top} \right)$$

$$+ \beta \operatorname{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_{2})^{\top}} \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(1;m_{1})} \mathbf{U}^{(1)} \right)$$

$$+ \gamma \operatorname{Tr} [\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}} \mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}] \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})} \right).$$

The obtained result is in exact agreement with the update rule (17) specified for the matrix $\mathbf{W}^{(1;m_1)}$. Collectively, this result and Lemma 1 establish that the proposed update rule guarantees the monotonic decrease of the original objective function.

In line with the strategy described in Proposition 7.2, two separate cases can be analyzed for the update rules of $\mathbf{W}^{(2;m_2)}$, $\mathbf{H}^{(1;m_1)}$ and $\mathbf{H}^{(2;m_2)}$. For each case, an auxiliary function is introduced to ensure the monotonic decrease of the original objective function. The cases are outlined as follows:

Case 1: Assuming that $\mathbf{W}^{(1;m_1)}$, $\mathbf{H}^{(1;m_1)}$, and $\mathbf{H}^{(2;m_2)}$ are fixed, the update rule (23) for $\mathbf{W}^{(2;m_2)}$ guarantees that the objective function in the minimization problem (12) is non-increasing. Under this scenario, the objective function with respect to $\mathbf{W}^{(2;m_2)}$ is expressed as

$$f(\mathbf{W}^{(2;m_2)}) = \frac{1}{2} \| \mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2$$

$$+ \frac{\beta}{2} \operatorname{Tr} [\mathbf{W}^{(2;m_2)} \mathbf{U}^{(2)} \mathbf{W}^{(2;m_2)^{\top}}] \operatorname{Tr} [\mathbf{W}^{(1;m_1)} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_1)^{\top}}]$$

$$+ \frac{\gamma}{4} \| \mathbf{W}^{(2;m_2)} \mathbf{W}^{(2;m_2)^{\top}} \otimes \mathbf{W}^{(1;m_1)} \mathbf{W}^{(1;m_1)^{\top}} - \mathbf{I}_{m_1 m_2} \|_F^2.$$

Next, by defining the function

$$g(w_{2}, w_{j_{2}, i_{2}}^{(2;m_{2})^{(t)}}) = \mathcal{B}(w_{j_{2}, i_{2}}^{(2;m_{2})^{(t)}}) + \dot{\mathcal{B}}(w_{j_{2}, i_{2}}^{(2;m_{2})^{(t)}})(w_{2} - w_{j_{2}, i_{2}}^{(2;m_{2})^{(t)}})$$

$$+ \left(\mathbf{H}^{(2;m_{2})^{\top}}\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})^{(t)}}\mathbf{X}_{(2)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})^{\top}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})\mathbf{X}_{(2)}^{\top}\right)$$

$$+ \beta \mathrm{Tr}[\mathbf{W}^{(1;m_{1})}\mathbf{U}^{(1)}\mathbf{W}^{(1;m_{1})^{\top}}]\mathbf{W}^{(2;m_{2})^{(t)}}\mathbf{U}^{(2)}$$

$$+ \gamma \mathrm{Tr}[\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}]\mathbf{W}^{(2;m_{2})^{(t)}}(\mathbf{W}^{(2;m_{2})^{(t)}})^{\top}\mathbf{W}^{(2;m_{2})^{(t)}})$$

$$\times \frac{(w_{2} - w_{j_{2}, i_{2}}^{(2;m_{2})^{(t)}})^{2}}{2 w_{j_{2}, i_{2}}^{(2;m_{2})^{(t)}}},$$

it can be demonstrated that $g(w_2, w_{j_2, i_2}^{(2;m_2)^{(t)}})$ serves as an auxiliary function for $\mathcal{B}(w_2)$, for $j_2=1,2,\ldots,m_2$, and $i_2=1,2,\ldots,I_2$. Note that $\mathcal{B}(w_2)$ represents the components of $f(\mathbf{W}^{(2;m_2)})$ associated with $w_{j_2,i_2}^{(2;m_2)}$ and takes the form

$$\mathcal{B}(w_2) = \mathcal{B}(w_{j_2,i_2}^{(2;m_2)^{(t)}}) + \dot{\mathcal{B}}(w_{j_2,i_2}^{(2;m_2)^{(t)}})(w_2 - w_{j_2,i_2}^{(2;m_2)^{(t)}}) + \frac{1}{2} \ddot{\mathcal{B}}(w_{j_2,i_2}^{(2;m_2)^{(t)}})(w_2 - w_{j_2,i_2}^{(2;m_2)^{(t)}})^2,$$

with

$$\dot{\mathcal{B}}(w_{2}) := \left(\frac{\partial f}{\partial \mathbf{W}^{(2;m_{2})}}\right)_{j_{2},i_{2}} = \left(-\mathbf{H}^{(2;m_{2})^{\top}}\mathbf{X}_{(2)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})\mathbf{X}_{(2)}^{\top} + \mathbf{H}^{(2;m_{2})^{\top}}\mathbf{H}^{(2;m_{2})}\mathbf{W}^{(2;m_{2})}\mathbf{X}_{(2)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})^{\top} \times (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})\mathbf{X}_{(2)}^{\top} + \beta \operatorname{Tr}[\mathbf{W}^{(1;m_{1})}\mathbf{U}^{(1)}\mathbf{W}^{(1;m_{1})^{\top}}]\mathbf{W}^{(2;m_{2})}\mathbf{U}^{(2)} + \gamma \operatorname{Tr}[\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}]\mathbf{W}^{(2;m_{2})} - \gamma \operatorname{Tr}[\mathbf{W}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})^{\top}}]\mathbf{W}^{(2;m_{2})}\right)_{i_{2},i_{2}},$$

and

$$\begin{split} \ddot{\mathcal{B}}(w_{2}) &:= \left(\frac{\partial^{2} f}{\partial \mathbf{W}^{(2;m_{2})^{2}}}\right)_{j_{2},i_{2}} = \left(\mathbf{H}^{(2;m_{2})^{\top}} \mathbf{H}^{(2;m_{2})}\right)_{j_{2},j_{2}} \\ &\times \left(\mathbf{X}_{(2)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})})^{\top} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top}\right)_{i_{2},i_{2}} \\ &+ \beta \text{Tr} [\mathbf{W}^{(1;m_{1})} \mathbf{U}^{(1)} \mathbf{W}^{(1;m_{1})^{\top}}] u_{j_{2},i_{2}}^{(2)} \\ &+ \gamma \text{Tr} [\mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}} \mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}}] \\ &\times \left(\left(\mathbf{W}^{(2;m_{2})^{\top}} \mathbf{W}^{(2;m_{2})}\right)_{i_{2},i_{2}} + w_{j_{2},i_{2}}^{(2,m_{2})^{2}} \\ &+ \left(\mathbf{W}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})^{\top}}\right)_{j_{2},j_{2}} \right) - \gamma \text{Tr} [\mathbf{W}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})^{\top}}]. \end{split}$$

Case 2: Assuming that $\mathbf{W}^{(1;m_1)}$, $\mathbf{W}^{(2;m_2)}$, and $\mathbf{H}^{(2;m_2)}$ are fixed, the update rule (19) for $\mathbf{H}^{(1;m_1)}$ guarantees that the objective function in the minimization problem (12) is non-increasing. Under this scenario, the objective function with respect to $\mathbf{H}^{(1;m_1)}$ is expressed as

$$f(\mathbf{H}^{(1;m_1)}) = \frac{1}{2} \| \mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2$$
$$+ \frac{\alpha}{2} \operatorname{Tr} [\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \operatorname{Tr} [\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}].$$

Next, by defining the function

$$\begin{split} g(h_1, h_{i_1, j_1}^{(1; m_1)^{(t)}}) &= \mathcal{B}(h_{i_1, j_1}^{(1; m_1)^{(t)}}) + \dot{\mathcal{B}}(h_{i_1, j_1}^{(1; m_1)^{(t)}})(h_1 - h_{i_1, j_1}^{(1; m_1)^{(t)}}) \\ &\quad + \left(\mathbf{H}^{(1; m_1)^{(t)}} \mathbf{W}^{(1; m_1)} \mathbf{X}_{(1)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2; m_2)} \mathbf{W}^{(2; m_2)})^{\top} \right. \\ &\quad \times \left. (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2; m_2)} \mathbf{W}^{(2; m_2)}) \mathbf{X}_{(1)}^{\top} W^{(2; m_2)^{\top}} \right. \\ &\quad + \alpha \operatorname{Tr}[\mathbf{H}^{(2; m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2; m_2)}] \mathbf{B}^{(1)} \mathbf{H}^{(1; m_1)} \right)_{i_1, j_1} \frac{(h_1 - h_{i_1, j_1}^{(1; m_1)^{(t)}})^2}{2 h_{i_1, j_1}^{(1; m_1)^{(t)}}}, \end{split}$$

it can be demonstrated that $g(h_1, h_{i_1, j_1}^{(1; m_1)^{(t)}})$ serves as an auxiliary function for $\mathcal{B}(h_1)$, for $i_1 = 1, \ldots, I_1$, and $j_2 = 1, \ldots, m_1$. Note that $\mathcal{B}(h_1)$ represents the components of $f(\mathbf{H}^{(1; m_1)})$ associated with $h_{i_1, j_1}^{(1; m_1)}$ and takes the form

$$\mathcal{B}(h_1) = \mathcal{B}(h_{i_1,j_1}^{(1;m_1)^{(t)}}) + \dot{\mathcal{B}}(h_{i_1,j_1}^{(1;m_1)^{(t)}})(h_1 - h_{i_1,j_1}^{(1;m_1)^{(t)}}) + \frac{1}{2}\ddot{\mathcal{B}}(h_{i_1,j_1}^{(1;m_1)^{(t)}})(h_1 - h_{i_1,j_1}^{(1;m_1)^{(t)}})^2,$$

with

$$\dot{\mathcal{B}}(h_1) := \left(\frac{\partial f}{\partial \mathbf{H}^{(1;m_1)}}\right)_{i_1,j_1} = \left(-\mathbf{X}_{(1)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})\mathbf{X}_{(1)}^{\top}\mathbf{W}^{(1;m_1)^{\top}} + \mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)}\mathbf{X}_{(1)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})^{\top} \times (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})\mathbf{X}_{(1)}^{\top}\mathbf{W}^{(1;m_1)^{\top}} + \alpha \text{Tr}[\mathbf{H}^{(2;m_2)^{\top}}\mathbf{L}^{(2)}\mathbf{H}^{(2;m_2)}]\mathbf{L}^{(1)}\mathbf{H}^{(1;m_1)}\right)_{i_1,i_1},$$

and

$$\ddot{\mathcal{B}}(h_1) := \left(\frac{\partial^2 f}{\partial \mathbf{H}^{(1;m_1)^2}}\right)_{i_1,j_1} = \left(\mathbf{W}^{(1;m_1)}\mathbf{X}_{(1)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})^{\top} \times (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})\mathbf{X}_{(1)}^{\top}\mathbf{W}^{(1;m_1)^{\top}}\right)_{j_1,j_1} \\
+ \alpha \text{Tr}[\mathbf{H}^{(2;m_2)^{\top}}\mathbf{L}^{(2)}\mathbf{H}^{(2;m_2)}]\ell_{i_1,i_1}^{(1)}.$$

Case 3: Assuming that $\mathbf{W}^{(1;m_1)}$, $\mathbf{W}^{(2;m_2)}$, and $\mathbf{H}^{(1;m_1)}$ are fixed, the update rule (25) for $\mathbf{H}^{(2;m_2)}$ guarantees that the objective function in the minimization problem (12) is non-increasing. Under this scenario, the objective function with respect to $\mathbf{H}^{(2;m_2)}$ is expressed as

$$f(\mathbf{H}^{(2;m_2)}) = \frac{1}{2} \| \mathcal{X} - \mathcal{X} \times_1 \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)} \times_2 \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)} \|_F^2$$
$$+ \frac{\alpha}{2} \operatorname{Tr}[\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \operatorname{Tr}[\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}].$$

Next, by defining the function

$$g(h_{2}, h_{i_{2}, j_{2}}^{(2;m_{2})^{(t)}}) = \mathcal{B}(h_{i_{2}, j_{2}}^{(2;m_{2})^{(t)}}) + \dot{\mathcal{B}}(h_{i_{2}, j_{2}}^{(2;m_{2})^{(t)}})(h_{1} - h_{i_{2}, j_{2}}^{(2;m_{2})^{(t)}})$$

$$+ \left(\mathbf{H}^{(2;m_{2})^{(t)}}\mathbf{W}^{(2;m_{2})}\mathbf{X}_{(2)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})^{\top}\right)$$

$$\times (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})}\mathbf{W}^{(1;m_{1})})\mathbf{X}_{(2)}^{\top}\mathbf{W}^{(2;m_{2})^{\top}}$$

$$+ \alpha \operatorname{Tr}[\mathbf{H}^{(1;m_{1})^{\top}}\mathbf{L}^{(1)}\mathbf{H}^{(1;m_{1})}]\mathbf{B}^{(2)}\mathbf{H}^{(2;m_{2})}\right)_{i_{2},i_{2}},$$

it can be demonstrated that $g(h_2,h_{i_2,j_2}^{(2;m_2)^{(t)}})$ serves as an auxiliary function for $\mathcal{B}(h_2)$, for $i_2=1,\ldots,I_2$, and $j_2=1,\ldots,m_2$. Note that $\mathcal{B}(h_2)$ represents the components of $f(\mathbf{H}^{(2;m_2)})$ associated with $h_{i_2,j_2}^{(2;m_2)}$ and takes the form

$$\mathcal{B}(h_2) = \mathcal{B}(h_{i_2,j_2}^{(2;m_2)^{(t)}}) + \dot{\mathcal{B}}(h_{i_2,j_2}^{(2;m_2)^{(t)}})(h_2 - h_{i_2,j_2}^{(2;m_2)^{(t)}}) + \frac{1}{2}\ddot{\mathcal{B}}(h_{i_2,j_2}^{(2;m_2)^{(t)}})(h_2 - h_{i_2,j_2}^{(2;m_2)^{(t)}})^2,$$

with

$$\dot{\mathcal{B}}(h_{2}) := \left(\frac{\partial f}{\partial \mathbf{H}^{(2;m_{2})}}\right)_{i_{2},j_{2}} = \left(-\mathbf{X}_{(2)}(\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_{2})^{\top}} + \mathbf{H}^{(2;m_{2})} \mathbf{W}^{(2;m_{2})} \mathbf{X}_{(2)} (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})})^{\top} \times (\mathbf{I}_{I_{3}} \otimes \mathbf{H}^{(1;m_{1})} \mathbf{W}^{(1;m_{1})}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_{2})^{\top}} + \alpha \text{Tr}[\mathbf{H}^{(1;m_{1})^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_{1})}] \mathbf{L}^{(2)} \mathbf{H}^{(2;m_{2})}\right)_{i_{2},j_{2}},$$

and

$$\ddot{\mathcal{B}}(h_2) := \left(\frac{\partial^2 f}{\partial \mathbf{H}^{(2;m_2)^2}}\right)_{i_2,j_2} = \left(\mathbf{W}^{(2;m_2)}\mathbf{X}_{(1)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)})^{\top} \times (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)})\mathbf{X}_{(2)}^{\top}\mathbf{W}^{(2;m_2)^{\top}}\right)_{j_2,j_2} \\
+ \alpha \text{Tr}[\mathbf{H}^{(1;m_1)^{\top}}\mathbf{L}^{(1)}\mathbf{H}^{(1;m_1)}]\ell_{i_2,i_2}^{(2)}.$$

7.5 COMPUTATIONAL COMPLEXITY

The purpose of this section is to evaluate the computational complexity of the suggested MSLFS method to offer a clear insight into its efficiency. Assessing the time complexity of each phase in Algorithm 1 allows for the calculation of the total computational expense. This evaluation will also emphasize the performance and scalability of the algorithm when managing large-scale applications. Initially, it is crucial to emphasize that for specific matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times r}$, $\mathbf{C} \in \mathbb{R}^{m \times n k}$, and $\mathbf{E} \in \mathbb{R}^{n \times n}$, the calculations for $\mathbf{A}\mathbf{B}$ and $\mathbf{C}(\mathbf{I}_k \otimes \mathbf{E})$ consist of 2mnr - mr and $2mn^2k - mnk$ arithmetic operations, respectively. It is important to note that the calculation $(\mathbf{I}_k \otimes \mathbf{E})$ requires no arithmetic operations since it is a diagonal matrix. Accordingly, the computational cost of updating the matrices $\mathbf{W}^{(1;m_1)}$, $\mathbf{H}^{(1;m_1)}$, $\mathbf{W}^{(2;m_2)}$ and $\mathbf{H}^{(2;m_2)}$ appears as follows:

1. The computational expense of updating the matrix $\mathbf{W}^{(1;m_1)}$ is Total flops($\mathbf{W}^{(1;m_1)}$)

$$\approx 6 \, m_1 I_1 I_2 I_3 + 4 \, I_1 I_2^2 I_3 + 2 \, m_1 I_2^2 I_3 + 6 \, m_2^2 I_2 + 6 \, m_2 I_2^2 + 8 \, m_1^2 I_1 + 2 \, m_1 I_1 + 2 \, m_2^3$$

$$= \mathcal{O}(m_1 I_1 I_2 I_3 + I_1 I_2^2 I_3) = \mathcal{O}(I_1 I_2 I_3 \max\{m_1, I_2\}).$$

2. The computational expense of updating the matrix $\mathbf{H}^{(1;m_1)}$ is Total flops($\mathbf{H}^{(1;m_1)}$)

$$\approx 8 m_1 I_1 I_2 I_3 + 6 I_1 I_2^2 I_3 + 8 m_2 I_2^2 + 4 m_1 I_1^2 + 2 m_2^2 I_2 + 4 m_1 I_1^2$$

= $\mathcal{O}(m_1 I_1 I_2 I_3 + I_1 I_2^2 I_3) = \mathcal{O}(I_1 I_2 I_3 \max\{m_1, I_2\}).$

3. The computational expense of updating the matrix $\mathbf{W}^{(2;m_2)}$ is Total flops($\mathbf{W}^{(2;m_2)}$)

$$\approx 6 m_2 I_1 I_2 I_3 + 4 I_2 I_1^2 I_3 + 2 m_2 I_1^2 I_3 + 6 m_1^2 I_1 + 6 m_1 I_1^2 + 8 m_2^2 I_2 + 2 m_2 I_2 + 2 m_1^3$$

= $\mathcal{O}(m_2 I_1 I_2 I_3 + I_1^2 I_2 I_3) = \mathcal{O}(I_1 I_2 I_3 \max\{I_1, m_2\}).$

4. The computational expense of updating the matrix $\mathbf{H}^{(2;m_2)}$ is Total flops($\mathbf{H}^{(2;m_2)}$)

$$\approx 8 m_2 I_1 I_2 I_3 + 6 I_2 I_1^2 I_3 + 8 m_1 I_1^2 + 4 m_2 I_2^2 + 2 m_1^2 I_1 + 4 m_2 I_2^2$$

= $\mathcal{O}(m_2 I_1 I_2 I_3 + I_2 I_1^2 I_3) = \mathcal{O}(I_1 I_2 I_3 \max\{m_2, I_1\}).$

To sum up, the computational expense of a single iteration of Algorithm 1 can be determined as follows:

Overall Total flops =
$$\mathcal{O}\bigg(I_1I_2I_3\big(\max\{m_1,I_2\}+\max\{m_2,I_1\}\big)\bigg)$$
.

7.6 MSLFS UPDATING RULES FOR A REAL-VALUED TENSOR DATA

To derive multiplicative updating rules when \mathcal{X} may be signed but all learned variables remain non-negative, we follow the same derivative computations as before and then apply elementwise positive/negative splitting to the matrix expressions that involve $\mathbf{X}_{(n)}$, where $n \in \{1,2\}$. For any real matrix \mathbf{M} we denote $\mathbf{M}_+ := \max(\mathbf{M},0)$ and $\mathbf{M}_- := \max(-\mathbf{M},0)$ (elementwise), so $\mathbf{M} = \mathbf{M}_+ - \mathbf{M}_-$. The multiplicative update rule for a non-negative variable \mathbf{Z} with gradient decomposed as $\nabla_{\mathbf{Z}}\mathcal{F} = \mathbf{G}^+ - \mathbf{G}^-$ (with $\mathbf{G}^\pm \geq 0$) is $\mathbf{Z} \leftarrow \mathbf{Z} \odot \mathbf{G}^- \oslash \mathbf{G}^+$. The gradients are fully developed in the previous section. Using the elementwise positive/negative splitting described above, the multiplicative updates (for non-negative factors while \mathcal{X} may be signed) are:

```
1091
                      \mathbf{W}^{(1;m_1)} = \mathbf{W}^{(1;m_1)} \odot [(\mathbf{H}^{(1;m_1)^\top} \mathbf{X}_{(1)} \mathbf{I}_{I_3} \otimes (\mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)}) \mathbf{X}_{(1)}^\top)_{-}
1092
                                              + \gamma \text{Tr}[\mathbf{W}^{(2;m_2)}\mathbf{W}^{(2;m_2)^\top}]\mathbf{W}^{(1;m_1)}]
1093
1094
                                              \oslash [(\mathbf{H}^{(1;m_1)^\top}\mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)}\mathbf{X}_{(1)}(\mathbf{I}_{I_3}\otimes\mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})^\top(\mathbf{I}_{I_3}\otimes\mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)})\mathbf{X}_{(1)}^\top)_+
1095
                                                    + \beta \text{Tr}[\mathbf{W}^{(2;m_2)}\mathbf{U}^{(2)}\mathbf{W}^{(2;m_2)^{\top}}]\mathbf{W}^{(1;m_1)}\mathbf{U}^{(1)}
1096
1097
                                                    + \gamma \text{Tr}[\mathbf{W}^{(2;m_2)}\mathbf{W}^{(2;m_2)^{\top}}\mathbf{W}^{(2;m_2)}\mathbf{W}^{(2;m_2)^{\top}}]\mathbf{W}^{(1;m_1)}\mathbf{W}^{(1;m_1)^{\top}}\mathbf{W}^{(1;m_1)}],
1099
                        \mathbf{H}^{(1;m_1)} = \mathbf{H}^{(1;m_1)} \odot [(\mathbf{X}_{(1)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2;m_2)} \mathbf{W}^{(2;m_2)}) \mathbf{X}_{(1)}^{\top} \mathbf{W}^{(1;m_1)^{\top}})_{-}
1100
                                              + \alpha \text{Tr}[\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \mathbf{A}^{(1)} \mathbf{H}_1]
1101
                                              \oslash \left[ (\mathbf{H}^{(1:m_1)}\mathbf{W}^{(1:m_1)}\mathbf{X}_{(1)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2:m_2)}\mathbf{W}^{(2:m_2))^\top} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(2:m_2)}\mathbf{W}^{(2:m_2)}) \mathbf{X}_{(1)}^\top \mathbf{W}^{(2:m_2)^\top} \right)_{+}
1102
1103
                                                    + \alpha \text{Tr}[\mathbf{H}^{(2;m_2)^{\top}} \mathbf{L}^{(2)} \mathbf{H}^{(2;m_2)}] \mathbf{B}^{(1)} \mathbf{H}^{(1;m_1)}],
1104
1105
                       \mathbf{W}^{(2;m_2)} = \mathbf{W}^{(2;m_2)} \odot [(\mathbf{H}^{(2;m_2)^\top} \mathbf{X}_{(2)} (\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)}) \mathbf{X}_{(2)}^\top)_{-}
1106
                                              + \gamma \text{Tr}[\mathbf{W}^{(1;m_1)}\mathbf{W}^{(1;m_1)^{\top}}]\mathbf{W}^{(2;m_2)}]
1107
1108
                                              \oslash [(\mathbf{H}^{(2;m_2)^\top}\mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)}\mathbf{X}_{(2)}(\mathbf{I}_{I_3}\otimes\mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)})^\top (\mathbf{I}_{I_3}\otimes\mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)})\mathbf{X}_{(2)}^\top)_{+}
1109
                                                    + \beta \text{Tr}[\mathbf{W}^{(1;m_1)}\mathbf{U}^{(1)}\mathbf{W}^{(1;m_1)^{\top}}]\mathbf{W}^{(2;m_2)}\mathbf{U}^{(2)}
1110
1111
                                                    + \gamma \text{Tr}[\mathbf{W}^{(1;m_1)}\mathbf{W}^{(1;m_1)^{\top}}\mathbf{W}^{(1;m_1)}\mathbf{W}^{(1;m_1)^{\top}}]\mathbf{W}^{(2;m_2)}\mathbf{W}^{(2;m_2)^{\top}}\mathbf{W}^{(2;m_2)}],
1112
                        \mathbf{H}^{(2;m_2)} = \mathbf{H}^{(2;m_2)} \odot [(\mathbf{X}_{(2)}(\mathbf{I}_{I_3} \otimes \mathbf{H}^{(1;m_1)} \mathbf{W}^{(1;m_1)}) \mathbf{X}_{(2)}^{\top} \mathbf{W}^{(2;m_2)^{\top}})_{-}
1113
1114
                                              + \alpha \text{Tr}[\mathbf{H}^{(1;m_1)^\top} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}] \mathbf{A}^{(2)} \mathbf{H}^{(2;m_2)}]
1115
                                              \oslash [(\mathbf{H}^{(2;m_2)}\mathbf{W}^{(2;m_2)}\mathbf{X}_{(2)}(\mathbf{I}_{I_3}\otimes\mathbf{H}^{(1;m_1)}\mathbf{W}^{(;m_1)})^\top (\mathbf{I}_{I_3}\otimes\mathbf{H}^{(1;m_1)}\mathbf{W}^{(1;m_1)})\mathbf{X}_{(2)}^\top\mathbf{W}^{(2;m_2)^\top})_+
1116
1117
                                                    + \alpha \text{Tr}[\mathbf{H}^{(1;m_1)^{\top}} \mathbf{L}^{(1)} \mathbf{H}^{(1;m_1)}] \mathbf{B}^{(2)} \mathbf{H}^{(2;m_2)}].
1118
```

8 ADDITIONAL EXPERIMENTAL RESULTS

8.1 Datasets

Table 5 summarizes the key statistics of the eight benchmark datasets used in our experiments, including the number of samples, feature dimensions, number of classes, and the range of selected features. These datasets together provide a comprehensive and diverse evaluation environment for assessing the proposed method across different domains, sample sizes, and feature complexities. COIL20 (Nene et al., 1996), ORL (Cai et al., 2010), and UMIST (Graham & Allinson, 1998) are classical image recognition benchmarks encompassing objects and human faces. COIL20 contains 20 objects imaged from multiple viewpoints, effectively testing robustness to pose variation. ORL consists of 40 subjects captured under relatively controlled conditions, whereas UMIST presents 20 subjects with more pronounced pose and illumination variations, creating a more challenging low-sample scenario. Pixraw10P and Orlraws10P (Li et al., 2017) are high-dimensional raw image subsets with limited samples, designed to evaluate the scalability of feature selection in situations where the number of features far exceeds the number of observations. Moving beyond traditional

object and face recognition, **FashionMNIST** (Xiao et al., 2017) serves as a modern drop-in replacement for the classic MNIST handwritten digit dataset, sharing the same grayscale 28×28 format but comprising clothing images with richer visual variability and finer inter-class distinctions, thus providing a more challenging benchmark while remaining compatible with MNIST's experimental protocols. In the biomedical domain, **BreastMNIST** and **OrganSMNIST** (Yang et al., 2021) focus on medical imaging tasks, with BreastMNIST providing a binary classification task based on breast ultrasound scans and OrganSMNIST involving multi-class organ recognition from MRI slices, thereby testing the applicability of the proposed approach to real-world medical scenarios. Collectively, these datasets span a wide range of sample sizes (from 100 to 1,440), feature dimensions (from 23×28 to 100×100), and class cardinalities (from 2 to 40), ensuring that the empirical evaluation thoroughly examines the method's robustness, scalability, and generalization ability across diverse, small-sample, high-dimensional, and domain-shifted settings.

Table 5: Detailed Statistics of the Eight Datasets.

Dataset	# of Samples	Feature Size	# of Classes	Range of Selected Features
COIL20	1,440	32×32	20	[50, 100,, 300]
ORL	400	32×32	40	[50, 100,, 300]
UMIST	575	23×28	20	[50, 100,, 300]
Pixraw10P	100	100×100	10	[50, 100,, 300]
Orlraws10P	100	92×112	10	[50, 100,, 300]
FashionMNIST	1,000	28×28	10	[50, 100,, 300]
BreastMNIST	546	28×28	2	[50, 100,, 300]
OrganSMNIST	500	28×28	11	[50, 100,, 300]

8.2 Comparison Models

This section summarizes the feature selection methods used for comparison, highlighting the core mechanism and strategy of each model to identify informative features while preserving relevant data structures.

- LS (He et al., 2005a): Assesses each feature individually based on how well it can maintain the local geometric structure of the data.
- UDFS (Yang et al., 2011): Selects the most informative features by performing both $\ell_{2,1}$ norm-based feature selection and local discriminative analysis at the same time.
- ILFS (Roffo et al., 2017): A probabilistic feature selection method that ranks features by considering all possible subsets while avoiding combinatorial complexity.
- **GRLTR** (Su et al., 2018): Combines low-rank tensor representation with local geometry preservation and $\ell_{2,1}$ norm-based feature selection.
- CAE (Balin et al., 2019): An end-to-end global feature selection approach that simultaneously trains a neural network to reconstruct the input data while selecting a representative subset of features.
- FSPCA (Tian et al., 2020): Simultaneously conducts feature selection and PCA by directly
 estimating the leading eigenvectors under row-sparsity constraints.
- CPUFS (Chen et al., 2023): Integrates a tensor-based linear classifier with graphregularized non-negative CP decomposition and pseudo-label regression.
- SPCAFS (Li et al., 2023): Applies a $\ell_{2,p}$ -norm sparsity regularization to the PCA projection matrix for feature selection.
- **GRSSLFS** (Tiwari et al., 2024): Selects high-variance basis features and integrates self-representation, subspace learning, and manifold regularization to enhance feature selection.
- SPDFS (Dong et al., 2025): Performs discriminative feature selection via ell_{2,0}-norm constrained sparse projection, combining fuzzy membership learning with globally and iteratively optimized projection strategies.

8.3 T-SNE VISUALIZATION ON ADDITIONAL DATASETS

In the main text, we presented t-SNE visualizations for the UMIST dataset. Here, we extend this analysis to Pixraw10P and Orlraws10P to further demonstrate the effectiveness of our unsupervised

tensor-based feature selector, MSLFS. Figure 5 shows the t-SNE embeddings of the original data and the embeddings obtained using the top 100, 200, and 300 features selected by MSLFS. As more informative features are included, intra-class samples become more tightly clustered while interclass samples separate more clearly, confirming that MSLFS effectively identifies discriminative features in an unsupervised manner.

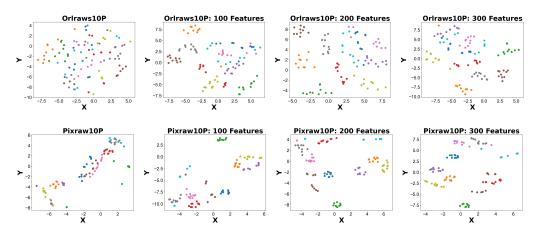


Figure 5: Visualization of t-SNE plots on the initial dataset and the dataset after applying the proposed model for feature reduction on the Orlraws10P and Pixraw10P datasets.

8.4 Customizing Feature Selection via Mode Combinations

To further evaluate the flexibility of MSLFS in distributing features across different tensor modes, we conducted an experiment on the Pixraw10P dataset by fixing the total number of selected features to 300 while varying the distribution of mode-1 and mode-2 slices. As illustrated in Figure 6, MSLFS can generate multiple valid configurations, such as 100×3 , 50×6 , or 10×30 , each corresponding to 300 intersection fibers. Across these different allocations, the selected slices capture meaningful vertical and horizontal structures. However, the results indicate that balanced selections across the two modes (e.g., 15×20 or 12×25) better preserve the overall inherent structure spanned by modes 1 and 2, while extreme allocations to a single mode tend to lose complementary information. This highlights that although MSLFS is flexible in how features are distributed, balanced configurations most effectively maintain both local and global structures.

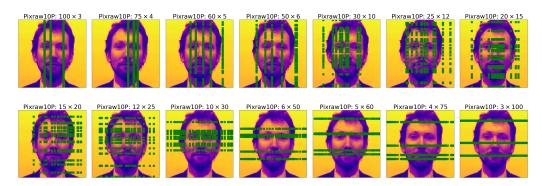


Figure 6: Visualization of mode-wise feature selection flexibility on Pixraw10P

8.5 Clustering on Selected Features

The experimental results with varying numbers of selected features further highlight the effectiveness of MSLFS. As shown in Figure 7, MSLFS is compared against 10 state-of-the-art models across eight benchmark datasets, where the performance curves illustrate both the absolute clustering accuracy and the stability of each method under different feature dimensions. Overall, MSLFS

consistently outperforms competing approaches, achieving the best or near-best results in terms of ACC and NMI across nearly all datasets. The improvements are especially notable on COIL20, ORL, and UMIST where MSLFS maintains clear superiority across varying feature subsets. Even on more challenging datasets such as BreastMNIST and OrganSMNIST, where existing methods often suffer from instability, MSLFS achieves significant margins, underscoring its robustness to data variability and imbalance. Furthermore, unlike other models that exhibit sharp fluctuations as the number of selected features changes, MSLFS demonstrates smooth and reliable performance trends, consistently producing discriminative feature subsets. This stability is largely attributed to its slice-based selection mechanism and higher-order graph regularization, which together preserve informative structures while effectively suppressing redundancy.

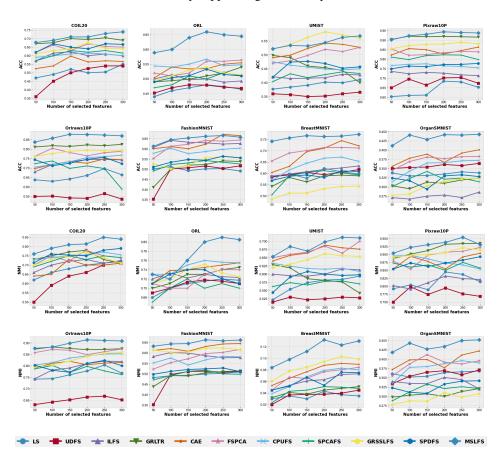


Figure 7: ACC and NMI curves of different feature selection methods on the eight datasets

8.6 SENSITIVITY ANALYSIS

To further investigate the influence of the regularization parameters α and β on the clustering performance of MSLFS, a sensitivity analysis is conducted. Figure 8 presents the heatmaps of NMI and ACC values across six datasets, including UMIST, Pixraw10P, Orlraws10P, ORL, OrganSM-NIST, and FashionMNIST. From Figure 8, it can be observed that the proposed method exhibits relatively stable behavior across a wide range of parameter values, though some dataset-specific trends emerge. For the UMIST dataset, both NMI and ACC remain stable with small fluctuations, and the best results are achieved when α lies within $\{10^1, 10^2, 10^3\}$ and β takes values around $\{10^1, 10^2, 10^3\}$. For the Pixraw10P dataset, MSLFS shows more sensitivity to β , with superior performance observed when $\alpha \in \{10^{-1}, 10^0, 10^1\}$ and β is set within $\{10^{-3}, 10^{-2}, 10^{-1}\}$. In the

Orlraws 10P dataset, MSLFS achieves consistently high NMI and ACC values, with optimal performance emerging when $\alpha \in \{10^{-3}, 10^{-2}, 10^3\}$ and $\beta \in \{10^{-2}, 10^{-1}, 10^0, 10^4\}$.

For the ORL dataset, the clustering performance is relatively insensitive to variations in β , while the most favorable results occur when α is chosen from $\{10^0, 10^1\}$. In the case of the OrganSMNIST dataset, both NMI and ACC show more noticeable fluctuations, but relatively better performance is achieved when $\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^4\}$ and β lies between $\{10^{-1}, 10^0, 10^1\}$. Finally, for the FashionMNIST dataset, the results indicate higher stability across parameter values, with the best performance obtained for $\alpha \in \{10^{-2}, 10^{-1}, 10^2, 10^2\}$ and $\beta \in \{10^{-1}, 10^0, 10^4\}$.

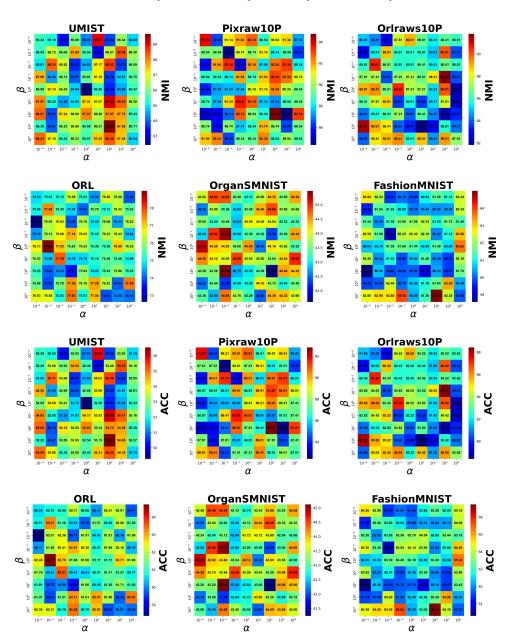


Figure 8: A comparison of the NMI and ACC scores obtained by MSLFS with different values of the parameters α , and β on six datasets.