

Multi-modal Differentiable Unsupervised Feature Selection

Junchen Yang¹ Ofir Lindenbaum² Yuval Kluger¹
Ariel Jaffe^{3†}

¹Yale University, USA;

²Bar-Ilan University, Israel;

³Hebrew University of Jerusalem, Israel

[†]Corresponding author. E-mail: ariel.jaffe@mail.huji.ac.il

Abstract

Multi-modal high throughput biological data presents a great scientific opportunity and a significant computational challenge. In multi-modal measurements, every sample is observed simultaneously by two or more sets of sensors. In such settings, many observed variables in both modalities are often nuisance and do not carry information about the phenomenon of interest. Here, we propose a multi-modal unsupervised feature selection framework: identifying informative variables based on coupled high-dimensional measurements. Our method is designed to identify features associated with two types of latent low-dimensional structures: (i) shared structures that govern the observations in both modalities and (ii) differential structures that appear in only one modality. To that end, we propose two Laplacian-based scoring operators. We incorporate the scores with differentiable gates that mask nuisance features and enhance the accuracy of the structure captured by the graph Laplacian. The performance of the new scheme is illustrated using synthetic and real datasets, including an extended biological application to single-cell multi-omics.

1 Introduction

In an effort to study biological systems, researchers are developing cutting-edge techniques that measure up to tens of thousands of variables at single-cell resolution. The complexity of such systems requires collecting multi-modal measurements to understand the interplay between different biological processes. Examples of such multi-modal measurements include SHARE-seq [1], DBiT-seq [2], CITE-seq [3], etc., which have provided biological insights and advancements in applications such as transcription factor characterization [4], cell type identification in human hippocampus [5], and immune cell profiling [6].

Multi-modal learning is a powerful tool widely used across multiple disciplines to extract latent information from high-dimensional measurements [7, 8]. Humans use complementary senses when attempting to “estimate” spoken words or sentences [9]. For example, lip movements can help us distinguish between two syllables that sound similar. The same intuition has inspired statisticians and machine learning researchers to develop learning techniques that exploit information captured simultaneously by complementary measurement devices.

Due to their applicability in multiple domains, there has been a growing interest in multi-modal approaches. Algorithms such as Contrastive Language–Image Pre-training (CLIP) [10], and Audioclip [11] have pushed the performance boundaries of machine learning for image, text, audio, analysis, and synthesis. The multi-modal data fusion task dates back to [12], which proposed the celebrated Canonical Correlation Analysis (CCA). CCA has many extensions [13, 14], and applications in diverse scientific domains [15, 16]. Despite their tremendous success, classical or advanced multi-modal schemes are often unsuitable for analyzing biological data. The large number of nuisance variables, which often exceeds the number of measurements, often causes correlation-based methods to overfit.

To attenuate the influence of nuisance or noisy features, several authors proposed unsupervised feature selection (UFS) schemes [17]. UFS seeks small subsets of informative variables in order to improve downstream analysis tasks, such as clustering or manifold learning. Empirical results demonstrate that informative features are often smooth with respect to some latent structure [18]. In practice, the smoothness of features can be evaluated based on how slowly they vary with respect to a graph [19]. Follow-up works exploited this idea to identify informative features [20, 21]. An alternative paradigm for UFS seeks subsets of features that can be used to reconstruct the entire data effectively [22].

While most fusion methods focus on extracting information shared between modalities, we propose a multi-modal UFS framework to identify features associated both with structures that appear in both modalities, and structures that are *modality-specific*, and appear in only one modality. To capture the shared structure, we construct a

symmetric shared graph Laplacian operator that enhances the shared geometry across modalities. We further propose differential graph operators that capture smooth structures that are not shared with the other modality. To perform multi-modal feature selection, we incorporate differentiable gates [23, 24] with the *shared* and *modality-specific* graph Laplacian scoring functions. This leads to a differentiable UFS scheme that attenuates the influence of nuisance features during training and computes a more accurate Laplacian matrix [25].

Our contributions are four folds: (i) Develop a *shared* and *modality-specific* Laplacian scoring operators. (ii) Motivate our operators using a product of manifolds model. (iii) develop and implement a differentiable framework for multi-modal UFS. (iv) Evaluate the merits and limitations of our approach with synthetic and real data and compare it to existing schemes.

2 Problem setting and preliminaries

We are given two data matrices $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$ whose rows contain n observations captured simultaneously in two modalities. The two sets of observations can be, for example, two arrays of sensors, cameras with different angles, etc. We are interested in processing modalities with bijective correspondences, which implies that there is a registration between the observations in both modalities.

Though the observations are high-dimensional, we assume that there are a small number of parameters governing the physical processes that underlies the data. These parameters can be continuous such as in a developmental process, or discrete - for example, when the observations can be characterized by clustering. However, the latent structure in both modalities may not be identical. For example, the two sets of observations may be generated by sets of sensors with different resolutions or sensitivity. For illustration, consider the observations shown in Fig. 1 (left). Both modalities follow a very similar tree structure. The bottom tree, however, has an additional bifurcating point that does not appear in the upper tree (green points).

Thus, we assume the latent parameters can be partitioned into two subsets. The first component denoted $\boldsymbol{\theta}_s$, captures the structures shared by both modalities. The second component, denoted $\boldsymbol{\theta}_x$ for modality \mathbf{X} , and $\boldsymbol{\theta}_y$ for modality \mathbf{Y} , captures the modality-specific structures that only appear in one set of observations. For example, the additional branch in the bottom tree (modality \mathbf{Y}) in Fig. 1 is governed by a parameter in $\boldsymbol{\theta}_y$. Thus, the observations \mathbf{X} and \mathbf{Y} are nonlinear transformations of $\boldsymbol{\theta}_s$, $\boldsymbol{\theta}_x$ and $\boldsymbol{\theta}_s$, $\boldsymbol{\theta}_y$, respectively.

Many biological data modalities are high dimensional and contain noisy features, which hinders the discovery of the underlying shared or modality-specific structures. Here, our goal is to identify groups of features associated with the shared structures $\boldsymbol{\theta}_s$

(e.g., the groups of features that are smooth on the shared bifurcated tree in Fig. 1) and groups of features associated with the modality-specific structures θ_x and θ_y (e.g., the features that are smooth with respect to the additional branch (θ_y) of modality \mathbf{Y} in Fig. 1). To achieve this goal, we compute two graphs that correspond to the two modalities. We use a spectral method to uncover the shared and graph-specific structures and apply a feature selection method to detect variables relevant to these structures. To better understand our approach, we first introduce some preliminaries about graph representation in Sec. 2.1, and discuss related work on feature selection in Sec. 2.2.

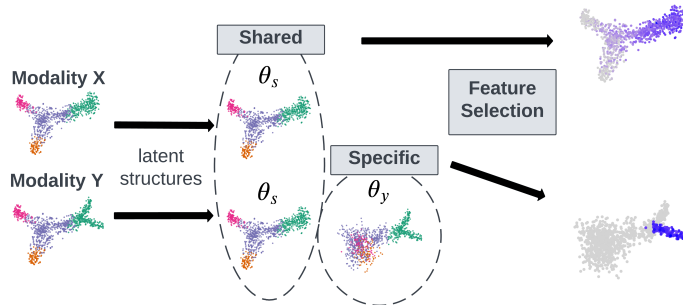


Figure 1: Overview of the goal: discovering features associated with shared and modality specific latent structures

2.1 The graph Laplacian and Laplacian score

A common assumption when analyzing high-dimensional datasets is that their structure lies on a low dimensional manifold in the high dimensional space [26, 27]. Methods for manifold learning are often based on a graph that captures the affinities between data points. Let $\mathbf{x}^{(i)}, \mathbf{y}^{(i)}$ denote the i -th observation in the \mathbf{X} and \mathbf{Y} modalities and let $\mathbf{K}_x, \mathbf{K}_y$ be, respectively, their affinity matrices whose elements are computed by the following Gaussian kernel functions.

$$(\mathbf{K}_x)_{i,j} = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\sigma_x^2}\right),$$

$$(\mathbf{K}_y)_{i,j} = \exp\left(-\frac{\|\mathbf{y}^{(i)} - \mathbf{y}^{(j)}\|^2}{2\sigma_y^2}\right),$$

where σ_x, σ_y are user-defined bandwidths that control the decay of each Gaussian kernel. Intuitively, the affinities decay exponentially with the distances between samples, thus capturing the local neighborhood structure in the high-dimensional space.

We compute the normalized Laplacian matrix by $\mathbf{L}_x = \mathbf{D}_x^{-\frac{1}{2}} \mathbf{K}_x \mathbf{D}_x^{-\frac{1}{2}}$, where \mathbf{D}_x is a diagonal matrix of row sums of \mathbf{K}_x . Similarly, \mathbf{L}_y is computed for modality \mathbf{Y} . An important property of the Laplacian matrix is that its eigenvectors corresponding to large eigenvalues reflect the underlying geometry of the data. The Laplacian eigenvectors are used for many applications including data embeddings [28], clustering [29], and feature selection [19]. For the latter, a popular metric for unsupervised identification of informative features is the Laplacian Score (LS) [19],

$$\mathbf{f}^T \mathbf{L}_x \mathbf{f} = \sum_{i=1}^n \lambda_i (\mathbf{f}^T \mathbf{u}_i)^2, \quad (1)$$

where $\mathbf{L}_x = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$ is the eigendecomposition of \mathbf{L}_x and \mathbf{f} is the normalized feature vector. Intuitively, when \mathbf{f} varies slowly with respect to the underlying structure of \mathbf{L}_x , it will have a significant component projected onto the subspace of its top eigenvectors, and a higher score.

2.2 Differentiable Unsupervised Feature Selection

A key limitation of the Laplacian score stems from its underlying assumption that the Laplacian matrix \mathbf{L}_x accurately reflects the latent structure of the data. This assumption, however, may not be valid in the presence of many noisy features. In such cases, the top eigenvectors of \mathbf{L}_x may be heavily influenced by noise and would not capture the underlying structure accurately. A recent work [25] addresses this problem by developing Differentiable Unsupervised Feature Selection (DUFS), a framework that estimates the Laplacian matrix while simultaneously selecting informative features using Laplacian scores. Specifically, DUFS computes a binary vector $\mathbf{s} \in \{0, 1\}^d$ that indicates which features are kept ($s_j = 1$) and which features are not ($s_j = 0$). Let $\Delta(\mathbf{s})$ denote a diagonal matrix with \mathbf{s} on the diagonal. At each iteration of DUFS, the Laplacian is computed based on $\tilde{\mathbf{X}} = \mathbf{X} \Delta(\mathbf{s})$, while simultaneously updating \mathbf{s} by optimizing over the following loss function.

$$\mathcal{L} = -\frac{1}{n} \text{Tr}[\tilde{\mathbf{X}}^T \mathbf{L}_{\tilde{\mathbf{x}}} \tilde{\mathbf{X}}] + \lambda \|\mathbf{s}\|_0, \quad (2)$$

where $\text{Tr}[\cdot]$ denotes the matrix trace. The first term equals the sum of Laplacian Scores across all features normalized by the total number of samples n in a training batch. The second term is a ℓ_0 regularizer that imposes sparsity to the number of selected features, with λ being a tunable parameter that controls the sparsity level. The output of DUFS is a list of a small number of selected features, and the Laplacian matrix $\mathbf{L}_{\tilde{\mathbf{x}}}$ learned from them.

However, due to the discrete nature of the ℓ_0 regularizer, the standard discrete indicator vector $\mathbf{s} \in \{0, 1\}^D$ will make objective in Eq. (2) not differentiable and

finding the optimal solution intractable. Following, [23], one can relax the ℓ_0 norm to a probabilistic differentiable counterpart, by replacing the binary indicator vector \mathbf{s} with a relaxed Bernoulli vector \mathbf{z} . Specifically, \mathbf{z} is a continuous Gaussian reparametrization of the discrete random variables, termed Stochastic Gates. It is defined for each feature i :

$$z_i = \max(0, \min(1, 0.5 + \mu_i + \epsilon_i)), \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad (3)$$

where μ_i is a learnable parameter, and σ is fixed throughout training. The loss function in Eq. (2) can now be reformulated as follows, which is the final objective of the DDFS:

$$\mathcal{L} = -\frac{1}{n} \text{Tr}[\tilde{\mathbf{X}}^T \mathbf{L}_{\tilde{x}} \tilde{\mathbf{X}}] + \lambda \|\mathbf{z}\|_0. \quad (4)$$

3 Method

We now derive our approach for unsupervised feature selection in multi-modal settings. Our method is designed to capture two types of features: (i) Features associated with latent structures that are *shared* between two modalities. (ii) Features associated with *differential latent structures*, that appear in only one modality. In Sec. 3.1 and 3.2, we derive two operators designed to capture shared and differential structures, respectively. To motivate our approach and illustrate the difference between shared and differential structures, we specifically address two examples: (i) shared and differential clusters and (ii) product of manifolds. We use the proposed operators in Sec. 3.3 to derive mmDDFS.

3.1 The shared structure operator

To motivate our approach, let us consider the artificial example illustrated in Fig. 2. The lower figure in the left panel shows the observations in modality \mathbf{Y} , which contains samples from a mixture of three distinct Gaussians. The upper figure shows modality \mathbf{X} , where one of the three clusters is partitioned again into three (less distinct) clusters.

It is instructive to study the *ideal setting* where we make the following assumptions: (i) The largest distance between two nodes within a cluster, denoted d_{within} is much smaller than the smallest distance between pairs of nodes of two clusters, denoted d_{between} . (ii) The bandwidth σ_x, σ_y is chosen such that $d_{\text{within}} \ll \sigma_x, \sigma_y \ll d_{\text{between}}$. In this setting, the three Gaussians constitute three main clusters, with no connections between pairs of nodes of different clusters and similar weights between pairs of nodes within clusters. Thus, the leading eigenvectors of \mathbf{L}_y span the subspace of the three *indicator vectors*. That is vectors that contain the square root of the degree of a node in a cluster and a zero value outside the cluster. See [29] and illustration in Fig. 2.

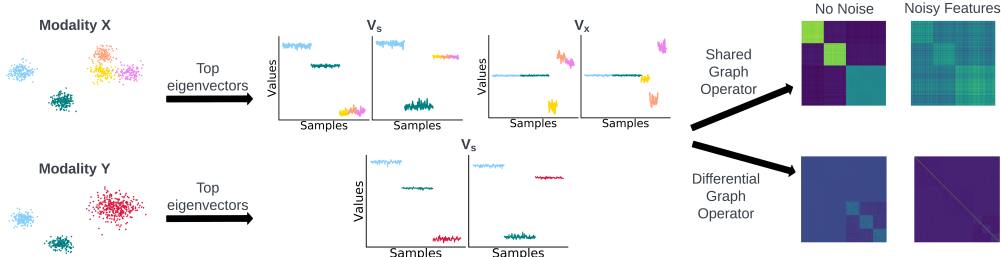


Figure 2: Visualization of the eigenvectors and the affinity matrix of the proposed operators on an artificial cluster example. Left: Visualization of the clusters. Middle: Leading eigenvectors of \mathbf{L}_x and \mathbf{L}_y . Right: Affinity matrices of the proposed shared graph operator (top) and the differential graph operator (bottom) with/without the presence of noisy features.

The matrix \mathbf{L}_x has two extra significant eigenvectors that span the separation of the third cluster, which appears only in \mathbf{X} . We denote by \mathbf{V}_s a matrix that contains the indicator vectors of the three partitions that appear in \mathbf{X} and \mathbf{Y} and by \mathbf{V}_x a matrix that contains the partitions that appear only in \mathbf{X} . In our ideal setting, the two Laplacian matrices $\mathbf{L}_x, \mathbf{L}_y$ are equal to

$$\mathbf{L}_x \approx \mathbf{V}_s \mathbf{V}_s^T + \mathbf{V}_x \mathbf{V}_x^T, \quad \mathbf{L}_y \approx \mathbf{V}_s \mathbf{V}_s^T. \quad (5)$$

To capture *shared* latent structures we compute the following shared operator $\mathbf{P}_{\text{shared}}$,

$$\mathbf{P}_{\text{shared}} = \mathbf{L}_x \mathbf{L}_y + \mathbf{L}_y \mathbf{L}_x. \quad (6)$$

For the cluster setting, the orthogonality between the matrices $\mathbf{V}_s, \mathbf{V}_x$ implies $\mathbf{P}_{\text{shared}} \approx 2\mathbf{V}_s \mathbf{V}_s^T$.

The symmetric product of the two Laplacians captures clusters that appear in both modalities while removing modality-specific clusters, see right panel of Fig. 2. We note that a similar operator to Eq. (6) is proposed in [30] for computing low-dimensional representations. Here, we combine our operator with DUFS to develop a multi-modal feature selection pipeline. We illustrate the usefulness of the shared operator for the product of manifold setting.

Product of manifolds. Let $\mathcal{M}_a, \mathcal{M}_b$ and \mathcal{M}_s be three low-dimensional manifolds embedded in \mathbb{R}^n , which are smooth transformations of three sets of latent variables $\boldsymbol{\theta}_a, \boldsymbol{\theta}_b$ and $\boldsymbol{\theta}_s$. To further motivate our approach, consider the case where modalities \mathbf{X} and \mathbf{Y} each contains observations from the products $\mathcal{M}_y, \mathcal{M}_x$ given by,

$$\mathcal{M}_y = \mathcal{M}_s \times \mathcal{M}_a, \quad \mathcal{M}_x = \mathcal{M}_s \times \mathcal{M}_b.$$

Note that the dependence on $\boldsymbol{\theta}_s$ is shared between $\mathcal{M}_x, \mathcal{M}_y$, while the dependence on $\boldsymbol{\theta}_a, \boldsymbol{\theta}_b$ is modality-specific.

In a product of manifolds $\mathcal{M}_x = \mathcal{M}_s \times \mathcal{M}_b$, every point $\boldsymbol{x} \in \mathcal{M}_x$ is associated with two points $\boldsymbol{x}_s \in \mathcal{M}_s$ and $\boldsymbol{x}_b \in \mathcal{M}_b$. Thus, we can define projection operators $\pi_b^x(\boldsymbol{x}), \pi_s^x(\boldsymbol{x})$ that map a point \boldsymbol{x} in \mathcal{M}_x to points in $\mathcal{M}_b, \mathcal{M}_s$, respectively. In addition, for every function $f^b : \mathcal{M}_b \rightarrow \mathbb{R}$ we define its extension to the product manifold \mathcal{M}_x by

$$(f^b \circ \pi_b^x)(\boldsymbol{x}) = f^b(\pi_b^x(\boldsymbol{x})).$$

An important property of a product \mathcal{M}_x is that the eigenfunctions $f_{l,m}^x$ of the Laplace Beltrami operator are equal to the pointwise product of the eigenfunctions of $\mathcal{M}_b, \mathcal{M}_s$, extended to \mathcal{M}_x .

$$f_{l,m}^x = (f_l^s \circ \pi_s^x)(f_m^b \circ \pi_b^x). \quad (7)$$

We refer to [31] for a detailed description of the properties of the product of manifolds. A simple example of a product of manifolds is a 2D rectangle area $(\boldsymbol{\theta}_s, \boldsymbol{\theta}_b) \in [0, l_s] \times [0, l_b]$. the projection π_s^x yields the first coordinate, while π_b^x yields the second. The eigenfunctions of the product with Neumann boundary conditions are equal to,

$$f_{l,m} = \cos(\pi l \boldsymbol{\theta}_s / l_s) \cos(\pi m \boldsymbol{\theta}_b / l_b). \quad (8)$$

Observations generated uniformly at random over the product of manifolds.

Here, we assume that the observations in the two modalities are generated by random and independent uniformly distributed samples over $\mathcal{M}_x, \mathcal{M}_y$. Let $\phi_{l,m}^x(\boldsymbol{x}_i), \phi_{l,k}^y(\boldsymbol{y}_i)$ denote the eigenvectors of $\boldsymbol{L}_x, \boldsymbol{L}_y$ evaluated at $\boldsymbol{x}_i, \boldsymbol{y}_i$ respectively. In the asymptotic regime where the number of points $n \rightarrow \infty$, the eigenvectors converge to the eigenfunctions as characterized in Eq. (7).

$$\begin{aligned} \phi_{l,m}^x(\boldsymbol{x}_i) &= \phi_l^s(\pi_s^x(\boldsymbol{x}_i)) \phi_m^b(\pi_b^x(\boldsymbol{x}_i)) \\ \phi_{l,k}^y(\boldsymbol{y}_i) &= \phi_l^s(\pi_s^y(\boldsymbol{y}_i)) \phi_k^a(\pi_a^y(\boldsymbol{y}_i)). \end{aligned} \quad (9)$$

Details about the definition and rate of convergence can be found, for example, in [32, 33], and reference therein. It is instructive to consider the ideal case, where due to their dependence on the independent projections π_b^x and π_a^x , the eigenvectors $\phi_{l,m}^x, \phi_{l,k}^y$ satisfy the following orthogonality property,

$$(\phi_{l,m}^x)^T \phi_{l',k}^y = \begin{cases} 1 & l = l', m = k = 0 \\ 0 & o.w. \end{cases} \quad (10)$$

It follows that the operator $\boldsymbol{P}_{\text{shared}}$ is equal to,

$$\boldsymbol{P}_{\text{shared}} = \boldsymbol{L}_x \boldsymbol{L}_y + \boldsymbol{L}_y \boldsymbol{L}_x = \sum_l (\phi_l^s \otimes \phi_0^a)(\phi_l^s \otimes \phi_0^b)^T, \quad (11)$$

where \otimes denotes the Hadamard product. The vectors ϕ_0^a, ϕ_0^b constitute the degree of the different observations and have little effect on the outcome. Thus, the leading eigenvectors of $\mathbf{P}_{\text{shared}}$ are associated with the shared component and not the differential components in the product of manifolds. Below, we illustrate this phenomenon with two examples.

Example 1: points in a 3D cube. Consider points generated uniformly at random over a 3D cube of dimensions $[0, l_s] \times [0, l_a] \times [0, l_b]$. Let $\mathbf{Y} \in \mathbb{R}^{n \times 2}$ constitute the first two coordinates of n independent observations, and let \mathbf{X} constitute the first and third coordinates. This is a simple case of a product of manifolds, where the shared variable θ_s is the first coordinate, while the modality-specific variables θ_a, θ_b are the second and third coordinates. Following Eq. (8), the eigenvectors of the graph Laplacian matrices $\mathbf{L}_x, \mathbf{L}_y$, evaluated at (θ_s, θ_b) and (θ_s, θ_a) converge to,

$$\begin{aligned}\phi_{lm}^x(\theta_s, \theta_b) &= \cos(\pi l \theta_s / l_s) \cos(\pi m \theta_b / l_b) \\ \phi_{lk}^y(\theta_s, \theta_a) &= \cos(\pi l \theta_s / l_s) \cos(\pi k \theta_a / l_a).\end{aligned}\tag{12}$$

The first row of Fig. 1 (Appendix A) shows a scatter plot of the points in \mathbf{X} (located according to the first two coordinates), colored by the values of the leading eigenvectors of \mathbf{L}_x . The second row shows the points in \mathbf{X} , but colored by the eigenvectors of $\mathbf{P}_{\text{shared}}$. As expected, all the eigenvectors of $\mathbf{P}_{\text{shared}}$ are functions of the shared coordinate θ_s .

Example 2: videos taken from different angles. Our second example is based on an experiment done in [34], where the two modalities constitute two videos of three dolls rotating at different angular speeds. The first camera (modality \mathbf{X}) captures the middle and left doll, while the second camera (modality \mathbf{Y}) captures the middle and right dolls (see Fig. 4a). Here, the shared variable θ_s is the angle of the middle doll captured by both modalities. The modality-specific variables θ_a, θ_b are the angles of the left and right dolls captured by each modality separately.

To illustrate Eq. (11) in this example, we first compute an approximation of the eigenvectors ϕ_i^s . To that end, we cropped each image in one of the videos such that only the middle doll (which appears in both modalities) is shown. One may think of this operation as a projection to the shared manifold. Next, we computed from the cropped images the leading eigenvectors ϕ_i^s of the Laplacian matrix. Fig. 2 (Appendix A) shows the leading three eigenvectors of $\mathbf{P}_{\text{shared}}$ as a function of $\phi_1^s, \phi_2^s, \phi_3^s$ as computed by the cropped images. The figure shows a linear dependency between the vectors, which implies that the shared operator retained only the shared component of the two modalities.

3.2 The Differential Graph Operators

We design two operators \mathbf{Q}_x and \mathbf{Q}_y to infer latent structures that are *modality specific* to \mathbf{X}, \mathbf{Y} respectively.

$$\mathbf{Q}_x = \tilde{\mathbf{L}}_y^{-1} \mathbf{L}_x \tilde{\mathbf{L}}_y^{-1}, \quad \mathbf{Q}_y = \tilde{\mathbf{L}}_x^{-1} \mathbf{L}_y \tilde{\mathbf{L}}_x^{-1}, \quad (13)$$

where $\tilde{\mathbf{L}}_x = \mathbf{L}_x + c\mathbf{I}$, $\tilde{\mathbf{L}}_y = \mathbf{L}_y + c\mathbf{I}$, and c is a regularization constant. We address the cluster example used for the shared operator to motivate the use of these operators.

Differential clusters. In the synthetic cluster example in Fig. 2, modality \mathbf{X} has three smaller clusters not observed in modality \mathbf{Y} . We show that one can detect the *differential clusters* of modality \mathbf{X} via the leading eigenvectors of \mathbf{Q}_x . By Eq. (5), we can approximate $\tilde{\mathbf{L}}_y$ via,

$$\tilde{\mathbf{L}}_y = (1 + c)\mathbf{V}_s \mathbf{V}_s^T + c\mathbf{V}_{\text{comp}} \mathbf{V}_{\text{comp}}^T, \quad (14)$$

where $\mathbf{V}_{\text{comp}} \in \mathbb{R}^{n \times (n-3)}$ contains, as columns, vectors that span the complementary subspace to \mathbf{V}_s . We write \mathbf{Q}_x as:

$$\mathbf{Q}_x = \tilde{\mathbf{L}}_y^{-1} \mathbf{L}_x \tilde{\mathbf{L}}_y^{-1} = c^{-2} \mathbf{V}_x \mathbf{V}_x^T + (1 + c)^{-2} \mathbf{V}_s \mathbf{V}_s^T. \quad (15)$$

The differential operator in Eq. (15) has two terms. The first spans the subspace corresponding to the differential structure \mathbf{V}_x , while the second spans the subspace of the shared structure \mathbf{V}_s . Since $c^{-2} > (1 + c)^{-2}$, it follows that the leading eigenvectors of \mathbf{Q}_x span the subspace of \mathbf{V}_x .

In theory, we can directly apply these operators to learn the structures. However, in many real-world applications, e.g., single-cell multi-omic technologies, both \mathbf{X} and \mathbf{Y} can be very noisy. In particular, abundant noisy features (e.g., genes) might dominate the data, and the top eigenvectors of \mathbf{L}_x and \mathbf{L}_y might not capture the underlying structure, which would be detrimental to the learning of $\mathbf{P}_{\text{shared}}$, \mathbf{Q}_x , and \mathbf{Q}_y . As shown in the affinity matrices on the right of Fig. 2, the structures are less clear when many noisy features are present. Therefore, it is necessary to have a feature selection framework that can effectively remove these noisy features in our multi-modal setting. With the aforementioned DUFFS feature selection framework as the foundation, we will show in the next section how we can incorporate it into our proposed operators in the multi-modal setting.

3.3 mmDUFFS

In this section, we describe our framework, termed multi-modal Differential Unsupervised Feature Selection (mmDUFFS)¹. We incorporate differentiable gates [25] with

¹Codes are available at <https://github.com/jcyang34/mmDUFFS>

loss functions based on the shared and differential operators, detailed in Sec. 3.1 and 3.2. Our goal is to compute an accurate shared graph operator ($\mathbf{P}_{\text{shared}}$ in Eq. (6)) and differential graph operators (\mathbf{Q}_x and \mathbf{Q}_y in Eq. (13)) while simultaneously selecting the informative features. Let $\mathbf{f}_x, \mathbf{f}_y$ denote a feature vector in \mathbf{X}, \mathbf{Y} , respectively. To quantify how noisy or informative the features are with respect to the shared structure, we replace the Laplacian \mathbf{L} in Eq. (1) with $\mathbf{P}_{\text{shared}}$, which yields the shared score $\mathbf{f}_x^T \mathbf{P}_{\text{shared}} \mathbf{f}_x$ and $\mathbf{f}_y^T \mathbf{P}_{\text{shared}} \mathbf{f}_y$. Similarly, $\mathbf{f}_x^T \mathbf{Q}_x \mathbf{f}_x$ and $\mathbf{f}_y^T \mathbf{Q}_y \mathbf{f}_y$ quantify the smoothness of these features with respect to the differential graph operators \mathbf{Q}_x and \mathbf{Q}_y . The rationale behind these generalized Laplacian Scores is similar to the original score. For instance, let $\mathbf{P}_{\text{shared}} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$ be the eigendecomposition of $\mathbf{P}_{\text{shared}}$. If \mathbf{f}_x varies slowly with respect to the underlying shared structure, it will have a larger component projected onto the subspace of $\mathbf{P}_{\text{shared}}$, thus leads to a higher score.

To learn features with high generalized Laplacian Scores and accurate graph operators, mmDUFS learns two sets of Stochastic Gates \mathbf{z}_x and \mathbf{z}_y that filter irrelevant features in each modality. Similar to DUFS [25], these stochastic gates multiply the data matrices \mathbf{X} and \mathbf{Y} to remove nuisance features, i.e., $\tilde{\mathbf{X}} = \mathbf{X} \Delta(\mathbf{z}_x)$ and $\tilde{\mathbf{Y}} = \mathbf{Y} \Delta(\mathbf{z}_y)$. At each iteration, the updated graph operators ($\tilde{\mathbf{P}}_{\text{shared}}, \tilde{\mathbf{Q}}_x, \tilde{\mathbf{Q}}_y$) are recomputed based on the gated inputs.

mmDUFS has two modes: (i) detecting shared structures using the shared graph operator $\tilde{\mathbf{P}}_{\text{shared}}$, and (ii) detecting modality-specific structures using the differential graph operators $\tilde{\mathbf{Q}}_x$, and $\tilde{\mathbf{Q}}_y$. To learn the shared structure and the corresponding features, we propose to optimize \mathbf{z}_x and \mathbf{z}_y by minimizing the following loss function:

$$\begin{aligned} \mathcal{L}_{\text{shared}} = & -\frac{1}{n} \text{Tr}[\tilde{\mathbf{X}}^T \tilde{\mathbf{P}}_{\text{shared}} \tilde{\mathbf{X}}] - \frac{1}{n} \text{Tr}[\tilde{\mathbf{Y}}^T \tilde{\mathbf{P}}_{\text{shared}} \tilde{\mathbf{Y}}] \\ & + \lambda_x \|\mathbf{z}_x\|_0 + \lambda_y \|\mathbf{z}_y\|_0, \end{aligned}$$

where the first two terms are the Shared Laplacian Scores for each modality, and the regularizers $\lambda_x \|\mathbf{z}_x\|_0$ and $\lambda_y \|\mathbf{z}_y\|_0$ control the number of selected features for each modality, with tunable parameters λ_x, λ_y that control the level of sparsity. In Appendix B.1, we suggest a procedure to tune these regularization parameters. Similarly, the loss functions $\mathcal{L}_x, \mathcal{L}_y$ are designed to detect features associated with structures that appear only in modality \mathbf{X}, \mathbf{Y} , respectively.

$$\begin{aligned} \mathcal{L}_x = & -\frac{1}{n} \text{Tr}[\tilde{\mathbf{X}}^T \mathbf{Q}_{\tilde{x}} \tilde{\mathbf{X}}] + \lambda_x \|\mathbf{z}_x\|_0, \\ \mathcal{L}_y = & -\frac{1}{n} \text{Tr}[\tilde{\mathbf{Y}}^T \mathbf{Q}_{\tilde{y}} \tilde{\mathbf{Y}}] + \lambda_y \|\mathbf{z}_y\|_0, \end{aligned} \tag{16}$$

where the first term in each loss is termed Differential Laplacian Scores. In the following section we show the usefulness of these score functions for detecting relevant features.

4 Results

We benchmark mmDUFs using synthetic and real multi-modal datasets. For discovering the shared structures and associated features, we compare mmDUFs with the shared operator to the following variants of kernel fusion-based methods previously proposed for dimensionality reduction: (1) Matrix Concatenation (MC), where the Laplacian is computed based on a concatenated matrix of the two modalities. (2) Multi-modal Kernel Sum (mmKS) [35], where the Laplacian is equal to $\mathbf{L}_x + \mathbf{L}_y$. (3) Multi-modal Kernel Product (mmKP) [36, 37], where the Laplacian is equal to $\mathbf{L}_x \mathbf{L}_y$.

For each baseline, the k features with the highest Laplacian Scores are selected. For the synthetic datasets, we set k to be the correct number of informative features. We evaluate the performance of different methods by the F1-score $F1 = TP / (TP + \frac{1}{2(FP+FN)})$, where TP is the number of informative features selected by each method, FP is the number of uninformative selected features, and FN is the number of missed informative features. For the rescaled MNIST and rotating doll examples, the informative features are set to the 25% pixels with the highest standard deviation.

4.1 Synthetic Examples

Rescaled MNIST. We designed a rescaled MNIST example with shared and modality-specific digits. We first randomly sample one image (28×28 pixels) of digits 0, 3, 8. Then, we rescale each digit randomly and independently 500 times resulting with 500 images of 0, 3, and 8. We concatenate pairs of 0 and 3 to create modality \mathbf{X} , and pairs of the same 3 and random 8 to create \mathbf{Y} , see example in Fig. 3a. Thus, this dataset consists of 500 samples and 28×56 pixels in each modality, with digit 3 shared between the modalities and digit 0 and 8 modality specific.

We apply mmDUFs with the shared operator to this example to select pixels corresponding to 3. The left column of Fig. 3b shows the pixels gate values from mmDUFs for modality \mathbf{X} (top) and \mathbf{Y} (bottom). We can see that selected pixels outline the shape of the digit 3 well. Table 1 compares the F1-score achieved by mmDUFs to three baselines. We can see that mmDUFs achieves a higher F1-score than all the baselines on both modalities, demonstrating its ability to identify informative features accurately.

Lastly, we apply mmDUFs with the differential operator to select modality-specific pixels. The right column of Fig. 3b shows the pixel gate values for both modality \mathbf{X} (top) and \mathbf{Y} (bottom). We can see that mmDUFs selects pixels that outline digits 0, 8 for modalities \mathbf{X} , \mathbf{Y} , respectively. Additionally, mmDUFs achieves F1-score 0.8059 and 0.8832 for \mathbf{X} and \mathbf{Y} , showcasing its effectiveness in identifying features contributing to the differential structures.

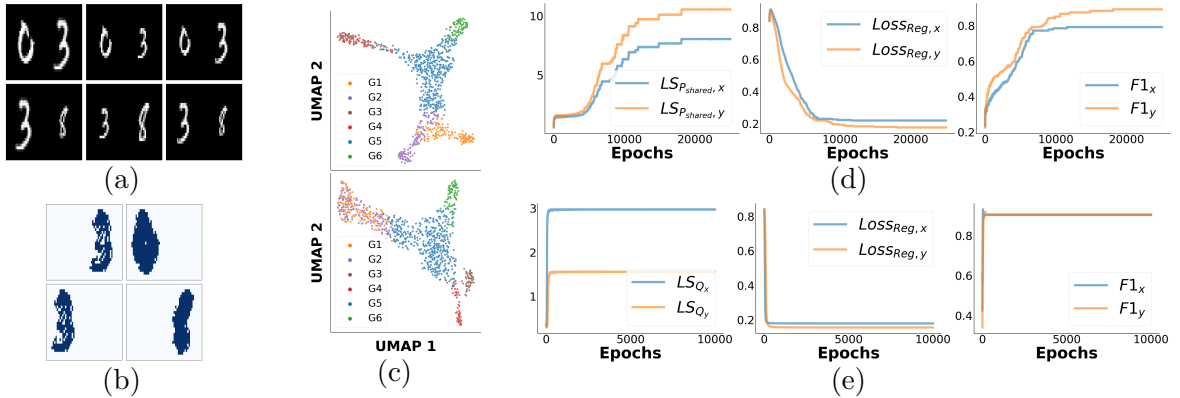


Figure 3: Left (a-b): Evaluation of the proposed approach on the rescaled MNIST dataset. (a): Random images from modality X (upper row) and modality Y (bottom row) in gray-scale. (b): Selected pixels (dark blue) for the shared operator (left column) and the differential operator (right column). Right (c-e): Synthetic developmental tree example. (c): UMAP embeddings of the tree using data from modality X (top) and modality Y (bottom). (d-e): Change of the Shared/Differential Laplacian Scores, regularization loss, and the F1-score of the selected features concerning the number of epochs (x-axis) for mmDUFFS with the shared operator (panel (c)) and the differential operator (panel (e)).

Dataset	Modality	MC	mmKS	mmKP	mmDUFFS
Rescaled MNIST	X	0.3547	0.5291	0.5291	0.7093
	Y	0.4826	0.6219	0.6219	0.8159
Synthetic Developmental Tree	X	0.6000	0.7800	0.8400	0.8800
	Y	0.7800	0.8000	0.8200	0.9000
Original Gaussian	X	0.5000	0.7333	1	1
	Y	0.5500	0.6500	0.9500	1
Gaussian + 10 Noisy Feats	X	0.5000	0.7333	1	1
	Y	0.5000	0.6500	0.9000	1
Gaussian + 30 Noisy Feats	X	0.4667	0.7000	0.9667	1
	Y	0.4500	0.5500	0.8500	1
Gaussian + 50 Noisy Feats	X	0.4000	0.6333	0.9333	0.9667
	Y	0.4000	0.5500	0.8000	0.8500

Table 1: Comparison of F1-score between different methods on the rescaled MNIST example, the synthetic tree example, and the Gaussian mixture example with different numbers of additive noisy features.

Synthetic Developmental Tree. Tree structures are ubiquitous throughout different biological processes and data modalities in single-cell biology [38, 39]. To understand the interplay of different mechanisms underlying the complex developmental process, it is vital to discover the genetic features that contribute to the tree structure shared across modalities and those that contribute to modality-specific structures.

We evaluate mmDUPS using a simulated developmental tree example generated via a tree simulator ². The original data has 1000 samples and 100 features. We divide the data into half, such that each modality has 50 informative features that contribute to the shared tree structure, as shown in the UMAP embeddings in Fig. 3c, where the samples in the tree are grouped into different branch groups (labeled G_1 to G_6). We then add 50 features drawn from negative binomial distributions to each modality to create differential branches, that are only observed in one modality. Specifically, branches G_1 and G_2 are bifurcated in modality \mathbf{X} (top UMAP embeddings) but are mixed in modality \mathbf{Y} (bottom UMAP embeddings), and G_3 and G_4 are bifurcated in modality \mathbf{Y} but are mixed in modality \mathbf{X} (see Supplementary section B.3 for further details). After log transformation and z-scoring the data, we concatenate 200 features drawn from $N(0, 1)$ to each modality as noisy features.

We apply our model with the shared and differential operators to recover the features that contribute to the overall tree structure and the set of features that contribute to the split branches, respectively. Fig. 3d shows the change, during training with the shared loss, in the Shared/Differential Laplacian Scores, the regularization loss, and the F1-score. Fig. 3e shows the same properties for the differential loss. Table 1 compares the F1-score of the selected features between different methods. Here as well, mmDUPS clearly outperforms the other methods.

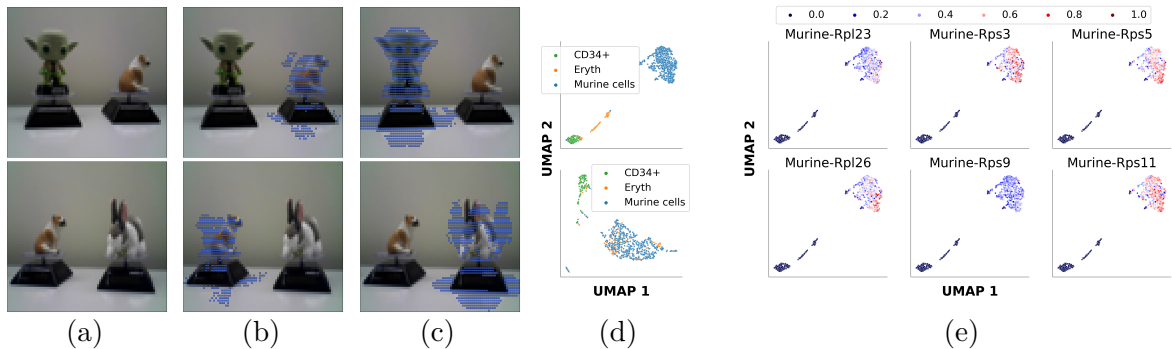


Figure 4: Left (a-c): Rotating dolls example. (a): Random images of the dolls from each video. (b-c): Selected pixels are marked in blue for mmDUPS with shared operator (b) and the differential operator (c). Right (d-e): CITE-seq data example. (d): UMAP embeddings using the RNA (top) and protein data (bottom), colored by cell type labels. (e): Similar UMAP embeddings colored by the expression level of several genes selected by mmDUPS with the differential operator.

²<https://github.com/dynverse/dyntoy>

Synthetic Gaussian Mixtures. We generated a multi-modal Gaussian mixture dataset, where \mathbf{X} and \mathbf{Y} each have 3 clusters. Two clusters are shared between modalities, and cluster 3 and 4 are specific to \mathbf{X} and \mathbf{Y} , respectively. Each cluster has a set of informative features drawn from a multivariate Gaussian, along with noisy features (see Appendix B.2 for details).

We first apply mmDUFFS to uncover the informative features of the shared clusters and the modality-specific clusters. In the figure of Supplementary section B.2, we plot the change of the average shared/differential Laplacian Scores across features, the regularization loss, and the F1-score of the selected features from mmDUFFS with respect to the number of epochs, where we can see that mmDUFFS gradually selects the correct features corresponding to high scores while sparsifying the number of features. To evaluate mmDUFFS’s feature selection capability in challenging regimes, we further inject 10, 30, and 50 noisy features into each modality and compare the F1-score of the selected features from different methods in each regime. As shown in Table 1, mmDUFFS consistently outperforms the baseline methods while maintaining accurate feature identification capability, demonstrating its robustness against noise.

4.2 Real Data

Rotating Dolls. We evaluate mmDUFFS’s performance on the rotating doll video dataset described in Sec. 3.1 in which 2 cameras capture 2 dolls from different angles (Fig. 4a). By treating each video frame as one sample (4050 in total) and the gray-scaled pixels as features, we aim to uncover pixels that correspond to the shared doll (the dog) and the modality-specific dolls (Yoda and rabbit).

For mmDUFFS with the shared operator, Fig. 4b shows selected pixels in both videos, as indicated by the blue dots. The shape of the dog is clearly delineated in both modalities. We further compute the F1-score of the selected pixels with respect to the underlying pixels that correspond to the dog. mmDUFFS achieves F1-score of 0.7158 and 0.8033 for the two modalities, whereas MC achieves 0.2390 and 0.3822, and mmKS and mmKP achieve 0.5452 and 0.6868. Fig. 4c shows the selected pixels of mmDUFFS with the differential operator in the two videos. In videos 1, mmDUFFS select mostly pixels corresponding to the Yoda (F1-score: 0.8861). For video 2, mmDUFFS select mostly pixels corresponding to the rabbit (F1-score: 0.7446).

CITE-seq Dataset. In single-cell biology, cell states are characterized by different features at different molecular levels. Identifying the contributing features is an open question crucial to understanding the underlying cell systems. We apply mmDUFFS to a CITE-seq dataset from [3], in which cells are profiled at both transcriptomic and proteomic levels measuring expressions of genes and protein markers, to identify the genes and proteins that characterize the cell states in the multi-modal setting.

In this data, a group of murine cells is spiked-in as controls to human cord blood mononuclear cells (CBMCs), and CITE-seq sequences the resulting cell system. Fig. 4d shows UMAP embeddings of the cells based on their RNA expression (top) and protein expression (bottom). From the full dataset, we analyzed 3 cell populations: murine cells (blue) and 2 CBMCs cell populations (Erythroids (orange) and CD34+ cells (green)). This dataset has 832 cells, with 500 top variable genes from modality 1 and 10 protein markers from modality 2. We can see that the murine cells are separable from the Erythroids in the RNA space but not in the proteomic space. To identify which gene markers contribute to the separation between cell groups, we apply mmDUFs with the differential operator to this data. We found that all the selected genes are murine genes that only express in the murine cells, as shown in Fig. 4e. This example demonstrates that mmDUFs can identify genetic markers contributing to the differential structures observed in single-cell multi-omic data.

5 Discussion

We present mmDUFs, a feature selection method that learns two novel graph operators that capture the *shared* and the *modality-specific* structures in multi-modal data, while simultaneously selecting the features that are informative for these structures. MmDUFs can operate on small batches which makes it scalable to large datasets. On the other hand, finding the optimal regularization parameters for mmDUFs on real data may be challenging, for which we suggest an automatic procedure in Appendix B.1. A second potential limitation is the $\mathcal{O}(n^3)$ computational complexity required to compute $\tilde{\mathbf{L}}$ (Eq. (13)). A possible solution is to reduce the complexity by computing a sparse Laplacian matrix.

Acknowledgements

The authors thank Amit Moscovich for the helpful discussions and feedback.

References

- [1] Sai Ma, Bing Zhang, Lindsay M LaFave, Andrew S Earl, Zachary Chiang, Yan Hu, Jiarui Ding, Alison Brack, Vinay K Kartha, Tristan Tay, et al. Chromatin potential identified by shared single-cell profiling of rna and chromatin. *Cell*, 183(4):1103–1116, 2020.

- [2] Yang Liu, Mingyu Yang, Yanxiang Deng, Graham Su, Archibald Enniful, Cindy C Guo, Toma Tebaldi, Di Zhang, Dongjoo Kim, Zhiliang Bai, et al. High-spatial-resolution multi-omics sequencing via deterministic barcoding in tissue. *Cell*, 183(6):1665–1681, 2020.
- [3] Marlon Stoeckius, Christoph Hafemeister, William Stephenson, Brian Houck-Loomis, Pratip K Chattopadhyay, Harold Swerdlow, Rahul Satija, and Peter Smibert. Simultaneous epitope and transcriptome measurement in single cells. *Nature methods*, 14(9):865–868, 2017.
- [4] Julia Joung, Sai Ma, Tristan Tay, Kathryn R Geiger-Schuller, Paul C Kirchgatterer, Vanessa K Verdine, Baolin Guo, Mario A Arias-Garcia, William E Allen, Ankita Singh, et al. A transcription factor atlas of directed differentiation. *Cell*, 186(1):209–229, 2023.
- [5] Yang Xiao, Graham Su, Yang Liu, Cheick A Sissoko, Yung-yu Huang, Adrienne N Santiago, Andrew J Dwork, Gorazd B Rosoklija, Underwood D Mark, Victoria Arango, et al. Spatially resolved transcriptomes in human hippocampus. *Biological Psychiatry*, 91(9):S18, 2022.
- [6] Noemie Leblay, Ranjan Maity, Elie Barakat, Sylvia McCulloch, Peter Duggan, Victor Jimenez-Zepeda, Nizar J Bahlis, and Paola Neri. Cite-seq profiling of t cells in multiple myeloma patients undergoing bcma targeting car-t or bites immunotherapy. *Blood*, 136:11–12, 2020.
- [7] Shiliang Sun. A survey of multi-view machine learning. *Neural computing and applications*, 23:2031–2038, 2013.
- [8] Xiaoqiang Yan, Shizhe Hu, Yiqiao Mao, Yangdong Ye, and Hui Yu. Deep multi-view learning methods: A review. *Neurocomputing*, 448:106–129, 2021.
- [9] Tommi Raij, Kimmo Uutela, and Riitta Hari. Audiovisual integration of letters in the human brain. *Neuron*, 28(2):617–625, 2000.
- [10] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [11] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 976–980. IEEE, 2022.

- [12] Harold Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.
- [13] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *International Conference on Machine Learning*, pages 1247–1255, 2013.
- [14] Ofir Lindenbaum, Moshe Salthov, Amir Averbuch, and Yuval Kluger. L0-sparse canonical correlation analysis. In *International Conference on Learning Representations*, 2022.
- [15] Harold Pimentel, Zhiyue Hu, and Haiyan Huang. Biclustering by sparse canonical correlation analysis. *Quantitative Biology*, 6(1):56–67, 2018.
- [16] Zhiwen Chen, Steven X Ding, Tao Peng, Chunhua Yang, and Weihua Gui. Fault detection for non-gaussian processes using generalized canonical correlation analysis and randomized algorithms. *IEEE Transactions on Industrial Electronics*, 65(2):1559–1567, 2017.
- [17] 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.
- [18] Alexandra Degeest, Michel Verleysen, and Benoît Frénay. Smoothness bias in relevance estimators for feature selection in regression. In *Artificial Intelligence Applications and Innovations: 14th IFIP WG 12.5 International Conference, AIAI 2018, Rhodes, Greece, May 25–27, 2018, Proceedings 14*, pages 285–294. Springer, 2018.
- [19] Xiaofei He, Deng Cai, and Partha Niyogi. Laplacian score for feature selection. *Advances in neural information processing systems*, 18, 2005.
- [20] Zheng Alan Zhao and Huan Liu. *Spectral feature selection for data mining*. Taylor & Francis, 2012.
- [21] Uri Shaham, Ofir Lindenbaum, Jonathan Svirsky, and Yuval Kluger. Deep unsupervised feature selection by discarding nuisance and correlated features. *Neural Networks*, 152:34–43, 2022.
- [22] Muhammed Fatih Balın, Abubakar Abid, and James Zou. Concrete autoencoders: Differentiable feature selection and reconstruction. In *International conference on machine learning*, pages 444–453. PMLR, 2019.

- [23] Yutaro Yamada, Ofir Lindenbaum, Sahand Negahban, and Yuval Kluger. Feature selection using stochastic gates. In *International Conference on Machine Learning*, pages 10648–10659. PMLR, 2020.
- [24] Junchen Yang, Ofir Lindenbaum, and Yuval Kluger. Locally sparse neural networks for tabular biomedical data. In *International Conference on Machine Learning*, pages 25123–25153. PMLR, 2022.
- [25] Ofir Lindenbaum, Uri Shaham, Erez Peterfreund, Jonathan Svirsky, Nicolas Casey, and Yuval Kluger. Differentiable unsupervised feature selection based on a gated laplacian. *Advances in Neural Information Processing Systems*, 34, 2021.
- [26] George C Linderman, Manas Rachh, Jeremy G Hoskins, Stefan Steinerberger, and Yuval Kluger. Fast interpolation-based t-sne for improved visualization of single-cell rna-seq data. *Nature methods*, 16(3):243–245, 2019.
- [27] Erez Peterfreund, Ofir Lindenbaum, Felix Dietrich, Tom Bertalan, Matan Gavish, Ioannis G Kevrekidis, and Ronald R Coifman. Local conformal autoencoder for standardized data coordinates. *Proceedings of the National Academy of Sciences*, 117(49):30918–30927, 2020.
- [28] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- [29] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [30] Tal Shnitzer, Mirela Ben-Chen, Leonidas Guibas, Ronen Talmon, and Hau-Tieng Wu. Recovering hidden components in multimodal data with composite diffusion operators. *SIAM Journal on Mathematics of Data Science*, 1(3):588–616, 2019.
- [31] Sharon Zhang, Amit Moscovich, and Amit Singer. Product manifold learning. In *International Conference on Artificial Intelligence and Statistics*, pages 3241–3249. PMLR, 2021.
- [32] Xiuyuan Cheng and Nan Wu. Eigen-convergence of gaussian kernelized graph laplacian by manifold heat interpolation. *Applied and Computational Harmonic Analysis*, 61:132–190, 2022.
- [33] Nicolás García Trillos, Moritz Gerlach, Matthias Hein, and Dejan Slepčev. Error estimates for spectral convergence of the graph laplacian on random geometric graphs toward the laplace–beltrami operator. *Foundations of Computational Mathematics*, 20(4):827–887, 2020.

- [34] Roy R Lederman and Ronen Talmon. Common manifold learning using alternating-diffusion. *submitted, Tech. Report YALEUIDCSITR1497*, 2014.
- [35] Dengyong Zhou and Christopher JC Burges. Spectral clustering and transductive learning with multiple views. In *Proceedings of the 24th international conference on Machine learning*, pages 1159–1166, 2007.
- [36] Ofir Lindenbaum, Arie Yeredor, and Moshe Salhov. Learning coupled embedding using multiview diffusion maps. In *Latent Variable Analysis and Signal Separation: 12th International Conference, LVA/ICA 2015, Liberec, Czech Republic, August 25-28, 2015, Proceedings 12*, pages 127–134. Springer, 2015.
- [37] Ofir Lindenbaum, Arie Yeredor, Moshe Salhov, and Amir Averbuch. Multi-view diffusion maps. *Information Fusion*, 55:127–149, 2020.
- [38] Mireya Plass, Jordi Solana, F Alexander Wolf, Salah Ayoub, Aristotelis Misios, Petar Glažar, Benedikt Obermayer, Fabian J Theis, Christine Kocks, and Nikolaus Rajewsky. Cell type atlas and lineage tree of a whole complex animal by single-cell transcriptomics. *Science*, 360(6391):eaaq1723, 2018.
- [39] Kai Zhang, James D Hocker, Michael Miller, Xiaomeng Hou, Joshua Chiou, Olivier B Poirion, Yunjiang Qiu, Yang E Li, Kyle J Gaulton, Allen Wang, et al. A single-cell atlas of chromatin accessibility in the human genome. *Cell*, 184(24):5985–6001, 2021.

A Additional Simulation Results

A.1 Points in a 3D cube.

The data consists of points in a 3D cube $[0, l_s] \times [0, l_a] \times [0, l_b]$. The modality \mathbf{X} includes the first two coordinates, and modality \mathbf{Y} includes the first and third, as explained in Sec. 3. The upper row in Figure A.1 shows the eigenvectors of \mathbf{L}_x . The eigenvectors change in both coordinates. The second row contains the eigenvectors of $\mathbf{P}_{\text{shared}}$. The leading eigenvectors change only with the first coordinate, as it is the only shared variable.

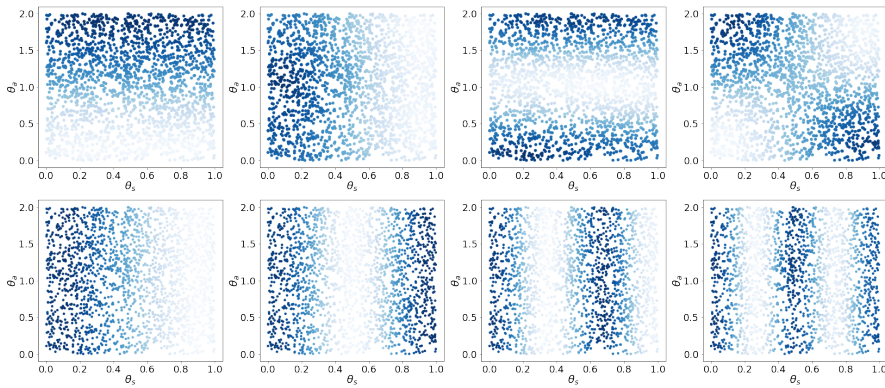


Figure A.1: Data consists of points sampled uniformly at random in a 3D cube. The upper row shows a scatter plot of the points, located according to the first two coordinates a, b and colored by the leading eigenvectors of \mathbf{L}_x , the Laplacian matrix of modality \mathbf{X} . The bottom row shows the leading eigenvectors of $\mathbf{P}_{\text{shared}}$, the product of Laplacians as defined in Eq. 6.

A.2 Rotating Dolls.

The two modalities include video frames taken simultaneously from two cameras, of three dolls rotating at different angular speeds. The first camera (modality \mathbf{X}) captures the left two dolls while the right camera (modality \mathbf{Y}) captures the right two dolls. Thus, the angle of the middle doll constitutes a shared variable θ_s . The angle of the left doll θ_x is modality \mathbf{X} -specific latent variable, and the angle of the right doll θ_y is modality \mathbf{Y} -specific latent variable.

From the left video, we cut the frames such that it includes only the middle doll (the shared component). From these images, we computed a graph Laplacian

matrix and its leading eigenvectors denoted ϕ_i^s . As explained in Sec. 3, we expect the eigenvectors of the shared operator, denoted v_i^s to be similar to ϕ_i^s , as both are associated with the latent variable θ_s . Figure A.2 shows v_i^s as a function of ϕ_i^s for $i = 1, 2, 3$. The three vectors are clearly highly correlated.

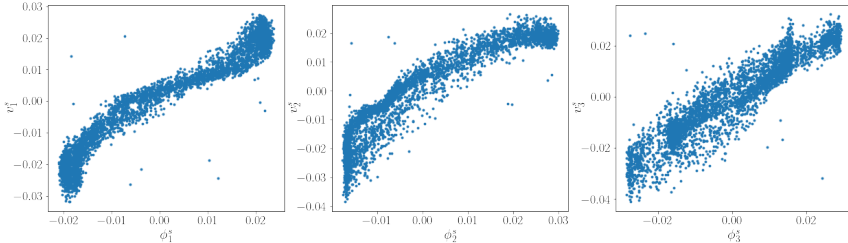


Figure A.2: The figure shows a scatter plot of v_i^s , the leading eigenvectors of $\mathbf{P}_{\text{shared}}$ as a function of ϕ_i^s , the estimated leading vectors of the shared component in the rotating doll dataset.

A.3 Synthetic Gaussian Mixtures.

Here we apply mmDUFs to uncover the informative features of the shared clusters and the modality-specific clusters. Fig. A.3b and Fig. A.3c show the change of the average Shared/Differential Laplacian Scores across features, the regularization loss, and the F1-score of the selected features from mmDUFs with respect to the number of epochs, where we can see that mmDUFs gradually selects the correct features corresponding to high scores while sparsifying the number of features.

B Experiment Details

In the following subsections, we provide additional experimental details required for the reproduction of the experiments provided in the main text. The CPU model used for the experiments is Intel(R) Xeon(R) Gold 6150 CPU @ 2.70GHz (72 cores total). The GPU model is NVIDIA GeForce RTX 2080 Ti.

Below in Table B.1 and B.2, we list the parameters we used on each experiment for mmDUFs with the shared operator and the differential operator. Parameter c is a regularization constant for mmDUFs with the differential operator, as mentioned in the main text. Parameter b is a scaling factor to the operators to balance between the Shared/Differential Laplacian Scores with respect to the regularization term. We used normalized Laplacian Matrix throughout the experiments except for the CITE-seq

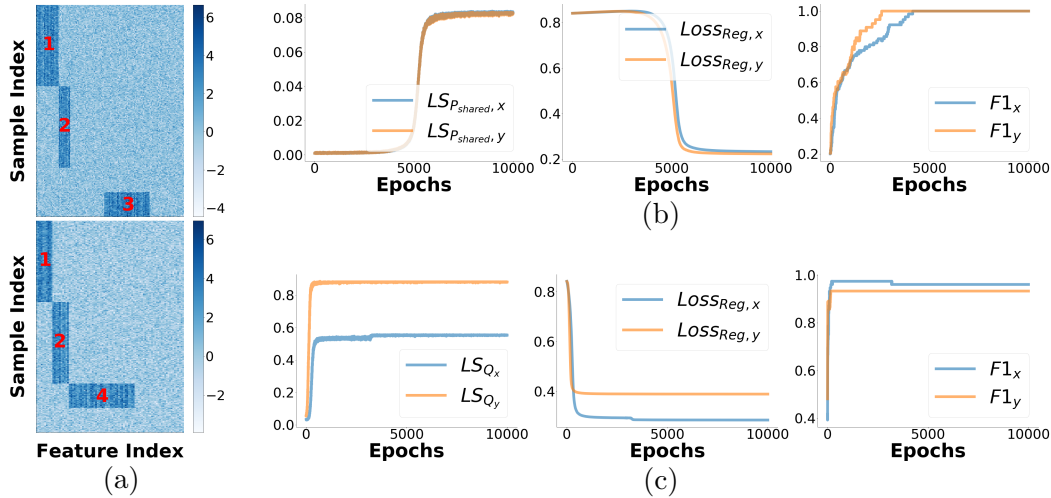


Figure A.3: Synthetic Gaussian mixture cluster example. (a): Data matrix of modality \mathbf{X} (top) and \mathbf{Y} (bottom). Rows are samples, and columns are features. Each modality has 3 clusters (labeled in red). Clusters 1 and 2 are shared between modalities, and cluster 3 and 4 are specific to each modality. (b): Change of the Shared Laplacian Scores, regularization loss, and the F1-score of the selected features concerning the number of epochs (x-axis) for mmDUFs with the shared operator. (c): Change of the Differential Laplacian Scores, regularization loss, and the F1-score of the selected features concerning the number of epochs (x-axis) for mmDUFs with the differential operator.

example where we found the performance was satisfactory with the un-normalized Laplacian Matrix.

Datasets	learning rate	epochs	λ_x	λ_y	b
Rescaled MNIST	2	10000	$1e-1$	$1e-1$	$1e2$
Synthetic Tree	2	25000	$1e-1$	$1e-1$	$1e3$
Gaussian Mixture	2	10000	$1e-4$	$1e-4$	1
Gaussian Mixture (10 Noisy Features)	2	20000	$1e-8$	$1e-6$	1
Gaussian Mixture (30 Noisy Features)	2	40000	$1e-4$	$1e-4$	1
Gaussian Mixture (50 Noisy Features)	2	10000	$1e-2$	$1e-3$	$1e2$
Rotating Dolls	2	10000	0.2	0.2	$1e3$

Table B.1: Parameters for mmDUFs with the shared operator across different datasets.

For the baseline methods, k features with the highest Laplacian Scores are selected. When evaluating the F1-score on the synthetic datasets, we set k to be the correct

Datasets	learning rate	epochs	λ_x	λ_y	c	b
Rescaled MNIST	1	10000	0.5	0.5	$1e-3$	$1e-4$
Synthetic Tree	2	10000	4	2	$1e-3$	$1e-3$
Gaussian Mixture	1	10000	0.4	0.4	$1e-1$	$1e-1$
Rotating Dolls	2	10000	2	2	3	$1e3$
CITE-seq	2	5000	3		2	1

Table B.2: Parameters for mmDUFFS with the differential operator across different datasets.

number of informative features. To make a fair comparison, we also let mmDUFFS select k features by sorting the raw gates (μ_d for feature d). For other datasets, we define selected features by mmDUFFS as features whose gates converged to 1 ($z_d = 1$ for feature d).

For the image datasets (rescaled MNIST, rotating dolls), we add small Gaussian noise drawn from $N(0, \sigma^2)$ to the pixels to stabilize feature selection of mmDUFFS. For the rescaled MNIST dataset, $\sigma = 0.1$ and we add noise to the non-informative pixels before standardizing the pixels via z-scoring. For the rotating dolls data, $\sigma = 5e-3$ and we add noise to all pixels before standardizing the pixels via z-scoring.

B.1 Tuning of the Regularization Parameter

mmDUFFS has tunable regularization parameters λ_x and λ_y that control the sparsity of the number of selected features. For synthetic datasets, one can tune these parameters to select features such that the selected number is close to the prescribed number s . However, it can still be time and resource-consuming to optimize these parameters. Also, for real data, one might not know how many features to select and what λ_x and λ_y to choose.

To alleviate this issue, we propose a "warm-up" procedure similar to [25] to optimize λ_x and λ_y . Specifically, we evaluate the mean Shared Laplacian Scores $S_{\text{shared}} = \frac{1}{2n}(\text{Tr}[\tilde{\mathbf{X}}^T \tilde{\mathbf{P}}_{\text{shared}} \tilde{\mathbf{X}}]/m + \text{Tr}[\tilde{\mathbf{Y}}^T \tilde{\mathbf{P}}_{\text{shared}} \tilde{\mathbf{Y}}]/d)$ and the mean Differential Laplacian Scores $S_x = \text{Tr}[\tilde{\mathbf{X}}^T \mathbf{Q}_{\tilde{x}} \tilde{\mathbf{X}}]/(d \times n)$, $S_y = \text{Tr}[\tilde{\mathbf{Y}}^T \mathbf{Q}_{\tilde{y}} \tilde{\mathbf{Y}}]/(m \times n)$ over a grid of λ_x and λ_y at the early stage of training (e.g., first 1000 epochs), and pick the parameters that maximize the Scores. Here n is the number of samples in the batch, and m and d are the number of selected features on each modality for real data or the number of pre-specified features for synthetic data.

To demonstrate this procedure, we use the synthetic Gaussian mixture dataset as the example, and we evaluate λ_x and λ_y over $\{1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-$

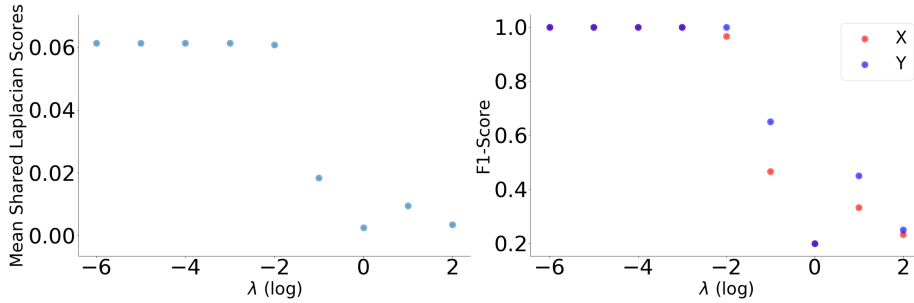


Figure B.4: Evaluation of the mean Shared Laplacian Scores (left) and the corresponding F1-scores (right) over a grid of λ s on the synthetic Gaussian mixture dataset. the y -axis shows the mean Shared Laplacian Scores (left) and F1-scores (right) whereas the x -axis shows the values of λ .

$1, 1, 1e1, 1e2\}$ using mmDUPS with the shared operator. For illustration purpose, we set $\lambda_x = \lambda_y$ Fig. B.4 shows the mean Shared Laplacian Scores over different λ values. We can see that $\{1e - 6, 1e - 5, 1e - 4, 1e - 3\}$ are the best candidates that give the highest Shared Laplacian Scores that also correspond to the highest F1-score.

B.2 Synthetic Gaussian Mixtures

We simulate 2 modalities \mathbf{X} and \mathbf{Y} , where modality \mathbf{X} has 260 samples with 130 features and modality \mathbf{Y} has 260 samples with 90 features. Both modalities have 3 clusters in the data (\mathbf{X} has cluster 1, 2, 3 and \mathbf{Y} has cluster 1, 2, 4, all labeled in red in Fig. A.3a), and each cluster has a set of informative features denoted as $\mathbf{f}_{x,i}$ and $\mathbf{f}_{y,i}$ ($i = 1, 2, 3, 4$) with length m_i ($i = 1, 2, 3, 4$). Each set of these informative features is drawn from $N(\boldsymbol{\mu}_i, \mathbf{I})$ independently for each sample, where $\boldsymbol{\mu}_i$ is a vector of length m_i drawn from $U(2, 4)$ and \mathbf{I} is an $m_i \times m_i$ identity matrix.

By design, cluster 1 and 2 are shared between modalities with $m_1 = 20$ and $m_2 = 10$ in modality \mathbf{X} , and $m_1 = 10$ and $m_2 = 10$ in modality \mathbf{Y} . On the other hand, cluster 3 is specific to modality \mathbf{X} with $m_3 = 40$, and cluster 4 is specific to modality \mathbf{Y} with $m_4 = 40$. The remaining features are considered noisy features and are drawn from $N(0, 1)$.

B.3 Synthetic Developmental Tree

We use `generate_data()` function from `dyntoy`³, a tree simulator package, to generate a dataset \mathbf{X}_0 with 1000 samples and 100 features. Specifically, the parameter

³<https://github.com/dynverse/dyntoy>

num_branchpoints is set to 1, *num_cells* is set to 1000, *num_features* is set to 100, *sample_mean_count* is set to 10, *sample_dispersion_count* is set to 50, *differentially_expressed_rate* is set to 4, and *dropout_probability_factor* is set to 0.

This step yields an initial data matrix $\mathbf{X}_0 \in \mathbb{R}^{1000 \times 100}$, and these 1000 samples are initially partitioned into 4 groups: G_1 and G_2 , G_3 and G_4 , G_5 , G_6 shown in Fig. 3c. For \mathbf{X}_0 , we further divide it into two halves, resulting in 2 data matrices $\mathbf{X} \in \mathbb{R}^{1000 \times 50}$ and $\mathbf{Y} \in \mathbb{R}^{1000 \times 50}$. We regard \mathbf{X} and \mathbf{Y} as 2 data modalities and these features as informative features contributing to the shared tree structure.

We further add 50 features to each modality that are drawn from negative binomial distributions to construct the differential structures between modalities. Specifically, for modality \mathbf{X} , the 50 features of G_1 are drawn from $NB(\mu = 4, \alpha = 0.1)$ where μ and α are the mean and dispersion parameter of the negative binomial distribution, whereas the 50 features of the other groups of samples are drawn from $NB(\mu = 20, \alpha = 0.1)$. Similarly, for modality \mathbf{Y} , the 50 features of G_3 are drawn from $NB(\mu = 4, \alpha = 0.1)$ while the 50 features of the other groups of samples are drawn from $NB(\mu = 20, \alpha = 0.1)$. Therefore, G_1 is bifurcated from G_2 and this structure is only observed in \mathbf{X} , and G_3 is bifurcated from G_4 and this structure is only observed in \mathbf{Y} .

Next, we row normalize each data matrix with a scaling factor $1e4$, and \log_{1p} transform the data. Then we standardize the features by z-scoring. At the end, we add 200 features drawn from $N(0, 1)$ to each modality as the noisy features.

B.4 CITE-seq

The human cord blood mononuclear cells (CBMCs) CITE-seq data was generated by [3], where the expression levels of both RNA and protein are measured for the same cells. We analyze 3 cell types: Erythroid cells, CD 34+ cells, and Murine cells. We row normalize each data matrix for both modalities. For the gene expression matrix (RNA), we filter the genes by standard deviation and keep the top 500 variable genes. Then for both matrices, we standardize the features by z-scoring.