

APPENDIX CONTENTS

- **Appendix A:** Dataset Details Page 15
- **Appendix B:** Implementation Details Page 16
- **Appendix C:** Additional Experiments Page 16
- **Appendix D:** Details on Khatri-Rao Clustering and Matrix Decomposition Page 19
- **Appendix E:** Proofs Page 22

Appendix A DATASET DETAILS

As stated in Section 9.1, in our experiments, we consider 13 synthetic and real-world datasets commonly used as benchmarks for clustering tasks. The datasets we consider reflect diverse data distributions, clustering structures and application domains. Next, we describe each dataset in detail.

- The MNIST dataset [15], available via PYTORCH [19] or CLUSTPY [17], is a collection of 28×28 (vectorized) grayscale images of handwritten digits. We draw a stratified subsample of 25000 images, ensuring that the original class proportions across the ten digit clusters are preserved. We rescale all pixel values by dividing by the maximum.
- The DOUBLE MNIST dataset is derived from MNIST by horizontally concatenating pairs of digit images. Each sample is a 28×56 (vectorized) grayscale image obtained by placing one 28×28 digit in the left position and another in the right position, and the label encodes the ordered pair of digits, yielding 100 clusters in total. For our experiments, we generate 10000 such composite images using the procedure described above, resulting in an approximately uniform distribution over all digit pairs. Unlike the MNIST dataset, the DOUBLE MNIST dataset by construction admits a natural clustering with a Khatri-Rao structure. We rescale all pixel values by dividing by the maximum.
- The HAR dataset, available via the UCI REPOSITORY [12] or CLUSTPY [17], consists of sensor readings collected from smartphone accelerometers and gyroscopes during human activity monitoring. There are 10299 data points. Each data point is represented as a multivariate feature vector of dimension 561 derived from raw time-series measurements, and the dataset includes 6 activity clusters such as walking, standing and sitting. We standardize each feature by subtracting its mean and dividing by its standard deviation.
- The OLIVETTI FACES dataset, available via SCIKIT-LEARN [20] or CLUSTPY [17], contains (vectorized) grayscale facial images of 40 individuals, with 10 images per subject captured under varying lighting conditions, facial expressions, and poses. Each image has resolution 64×64 pixels. We standardize each feature by subtracting its mean and dividing by its standard deviation.
- The CMU FACES dataset, available via the UCI REPOSITORY [12] or CLUSTPY [17], contains grayscale facial images of 20 persons varying their pose (up, straight, left and right) and expression (neutral, happy, sad, angry), and shown with and without sunglasses. In total, there are 624 images. The

original images with a resolution of 30×32 pixels are vectorized and each feature is standardized by subtracting its mean and dividing by its standard deviation.

- The SYMBOLS dataset, available via CLUSTPY⁶ [17], consists of vectorized handwritten symbols. Each sample corresponds to a time series obtained from the drawing trajectory of a symbol. The symbols are drawn by 13 different individuals. The dataset contains a total of 1020 data points, each with 398 measurements. The data points are organized into 6 natural clusters. We standardize each feature by subtracting its mean and dividing by its standard deviation.
- The STICKFIGURES dataset [9], available via CLUSTPY [17], consists of 900 synthetic (vectorized) silhouette images of human stick figures generated under varying poses. Each image has a resolution of 20×20 pixels. Examples of such images are given in Figure 1. Each image is provided at a fixed resolution and captures a simplified body configuration defined by joint positions and limb orientations, resulting in 9 distinct pose-based clusters. We rescale all pixel values by dividing by the maximum.
- The OPTDIGITS dataset, available via the UCI REPOSITORY [12] or CLUSTPY [17], contains 5620 (vectorized) handwritten digit images represented as 8×8 grayscale pixel grids. We standardize each feature by subtracting its mean and dividing by its standard deviation.
- The CLASSIFICATION dataset, sourced by SCIKIT-LEARN [20], is a synthetic dataset. While it was originally designed for benchmarking classifiers, it is also useful for evaluating clustering algorithms which simply discard class labels until the evaluation of the clustering results. Each class for classification corresponds to a cluster. Except in the experiments where we vary the number of data points, features, or clusters, the dataset consists of 5000 samples organized into 100 clusters. The data are generated with 10 informative features, with no redundant or repeated features, ensuring that all dimensions contribute meaningfully to class separability. We standardize each feature by subtracting its mean and dividing by its standard deviation.
- The CHAMELEON dataset, available via the CLUSTBENCH Benchmark Suite [8], consists of 10000 two-dimensional point clouds exhibiting complex, nonconvex cluster shapes with varying densities [11]. The dataset corresponds to the configuration shown at the bottom of Figure 4, augmented with a substantial proportion of uniformly distributed data points that do not belong to any natural cluster.
- The SOYBEAN LARGE dataset, available via the UCI REPOSITORY [12] or CLUSTPY [17], is a dataset of categorical features. It consists of 562 plant samples belonging to one of 15 classes. For each plant, there are 35 observed categorical attributes describing the plant. After encoding categorical attributes numerically, we standardize each feature by subtracting its mean and dividing by its standard deviation.
- The BLOBS dataset, sourced by SCIKIT-LEARN [20], is a synthetic dataset specifically designed for controlled evaluation of clustering algorithms that consists of data points grouped

⁶The dataset is also available at: <https://www.timeseriesclassification.com/index.php>

around isotropic Gaussian clusters. Each cluster has standard deviation 1. The dataset consists of 5000 2-dimensional data points, arranged in 100 clusters, except in the experiments where we vary the number of data points, features and clusters. We standardize each feature by subtracting its mean and dividing by its standard deviation.

- The R15 dataset, available via the CLUSTBENCH benchmark suite [8], consists of 600 two-dimensional data points forming 15 Gaussian clusters. The cluster centroids in R15 are not arranged on a regular grid but exhibit non-uniform inter-cluster distances. We standardize each feature by subtracting its mean and dividing by its standard deviation.

Appendix B IMPLEMENTATION DETAILS

In this section, we discuss implementation details for all the algorithms considered in the experiments presented in Section 9.

Implementation details of the naïve approach to Khatri-Rao clustering. For the naïve Khatri-Rao clustering baseline used in our experiments, we first run the SCIKIT-LEARN [20] implementation of standard k -MEANS to extract the desired number of cluster centroids. For instance, if the goal is to extract two sets of h_1 and h_2 protocentroids, k -MEANS retrieves $h_1 h_2$ clusters. Then, using a coordinate-descent procedure implemented in Python with closed-form updates, we decompose a set of centroids into two smaller sets of protocentroids whose Khatri-Rao product approximates the original centroids. The coordinate-descent procedure alternates between updating the protocentroids of the first and second set, holding the other set fixed. At each step, we use the closed-form updates illustrated in Section 5 (Equation (8)). The procedure stops either when a maximum number of iterations is reached (5000 by default) or when the total sum of squared differences between the initial centroids and the Khatri-Rao aggregation of the protocentroids becomes smaller than a user-specified threshold (10^{-4} by default). Upon termination of the described coordinate-descent procedure, we have the output sets of protocentroids. The corresponding centroids are readily obtained by aggregating protocentroids. To conclude, we assign each data point to the closest centroid. The implementation is available online in our code repository ⁷.

Implementation details of standard k -MEANS. We rely on the well-established SCIKIT-LEARN implementation of standard k -MEANS. The two k -MEANS baselines (namely k -MEANS with $h_1 + h_2$ centroids and $h_1 h_2$ centroids) are obtained by simply specifying the number of centroids $h_1 + h_2$ and $h_1 h_2$ as input. In the scalability experiments, instead, to ensure a fair comparison, we use an implementation of k -MEANS which mirrors the implementation of KHATRI-RAO- k -MEANS described next.

Implementation details of KHATRI-RAO- k -MEANS. Our experiments rely on a simple Python implementation of KHATRI-RAO- k -MEANS that is purely built on NUMPY ⁸, taking advantage of vectorized operations for efficiency. Such implementation is available online ⁹.

For initialization, by default we sample random data points as the initial protocentroids (as in Algorithm 1). Alternatively, we can

adopt the strategy inspired by k -MEANS++ described in Section 6, which selects representative data points based on squared-distance criteria. Our implementation of this strategy either deterministically chooses the data point farthest from the previously selected centroids, or samples data points with probability proportional to their distance from those centroids. After initialization, we iteratively compute assignments, protocentroid and corresponding centroids. Thanks to the closed-form updates introduced in Section 6, the updates of protocentroids (and centroids) are implemented in a fully-vectorized manner. At each iteration we monitor convergence by tracking the movement of all reconstructed centroids; the algorithm terminates once this movement falls below a user-specified threshold or a maximum number of iterations is reached.

During the execution of the algorithm, in case empty clusters arise, they are handled by reinitializing the corresponding protocentroid to a random data point.

Our implementation is deliberately simple and easy to follow. In the future, more optimized or parallelized implementations could be developed for larger-scale settings.

KHATRI-RAO- k -MEANS admits both a time-efficient and a memory-efficient implementation. Algorithm 1 presents the memory-efficient implementation that avoids storing the full set of centroids by computing them on the fly from the stored set of protocentroids. This approach can reduce memory requirements, particularly when the number of clusters is large since memory requirements only grow additively with the total number of protocentroids instead of multiplicatively. For example, $h_1 h_2$ clusters demand storing only $h_1 + h_2$ protocentroids instead of $h_1 h_2$ centroids. However, computing centroids on the fly incurs a runtime overhead. A time-efficient implementation is obtained by computing centroids once and storing them.

Implementation details of standard deep clustering algorithms. For DKM and IDEC, we rely on the off-the-shelf PYTORCH-based implementations provided by the CLUSTPY library ¹⁰. To find initial centroids, DKM and IDEC use the SCIKIT-LEARN implementation of standard k -MEANS.

Implementation details of Khatri-Rao deep clustering algorithms. For the implementation of Khatri-Rao deep clustering algorithms, we build the implementation of KHATRI-RAO DKM and KHATRI-RAO IDEC on top of the CLUSTPY implementations of the corresponding standard deep clustering algorithms.

Extending the CLUSTPY implementation of a standard deep clustering algorithm to the Khatri-Rao clustering paradigm is straightforward. During initialization, we rely on our implementation of KHATRI-RAO- k -MEANS to obtain initial protocentroids. We then reparameterize the centroids to satisfy the Khatri-Rao structure and adjust the autoencoder parameters to conform to the Hadamard-decomposition reparameterization. Our implementation of Khatri-Rao deep clustering algorithms is available online ¹¹.

Appendix C ADDITIONAL EXPERIMENTS

This section reports additional experimental results that complement the evaluation presented in Section 9. First, we empirically

⁷<https://github.com/maciap/KhatriRaoClustering/blob/main/scripts/KRkmeansExperimentsLib.py>

⁸<https://numpy.org>

⁹<https://github.com/maciap/KhatriRaoClustering/tree/main/KhatriRaoKMeans>

¹⁰<https://github.com/collinleiber/ClustPy>

¹¹<https://github.com/maciap/KhatriRaoClustering/tree/main/KhatriRaoDeepClustering>

assess the scalability of KHATRI-RAO- k -MEANS by measuring runtime and peak memory usage as the clustering tasks grow in size. Second, we study the effect of the number of centroids on the clustering performance of KHATRI-RAO- k -MEANS and the corresponding baselines. Finally, we turn to deep clustering and evaluate how compressing the autoencoder via the Hadamard decomposition affects clustering performance and training stability. Unless specified otherwise, all experimental settings are the same as for the experiments presented in Section 9.

Scalability. The goal of Khatri-Rao clustering is to extract accurate data summaries that are more succinct than those extracted by standard clustering methods. Our experiments show that Khatri-Rao clustering can yield more succinct clustering-based summaries while maintaining comparable accuracy. However, it is also important to ensure that the advantages provided by Khatri-Rao clustering are not undermined by significantly worse scalability. Section 6 discusses the theoretical asymptotic time and space complexity of KHATRI-RAO- k -MEANS. According to the discussion, KHATRI-RAO- k -MEANS has the same asymptotic time complexity as k -MEANS. In addition, KHATRI-RAO- k -MEANS can reduce the memory usage of k -MEANS when the number of clusters and hence centroids to be represented grows. To complement the analysis of time and space complexity, we conduct an empirical scalability analysis.

The results of the empirical scalability analysis are summarized in Figures 10 and 11. The figures show runtime (in seconds) for a single execution of an algorithm and peak memory usage (in Mebibytes) as the number of data points, features and centroids increase for the BLOBS and CLASSIFICATION datasets. In particular, we vary the number of data points in {4000, 8000, 12000, 16000} with 100 clusters and 100 features, we vary the number of features in {10000, 15000, 20000, 25000} with 1000 data points and 100 clusters, and we vary the number of clusters in {2500, 5000, 7500, 10000} with 20000 data points and 100 features.

The results suggest that KHATRI-RAO- k -MEANS clustering introduces a runtime overhead compared to standard k -MEANS, but this overhead does not significantly increase as the problem size increases. Instead, the overhead remains nearly constant as the number of data points, features and centroids increase, in line with the time-complexity discussion in Section 6.

The two-phase naïve approach to Khatri-Rao k -Means is usually the slowest algorithm, except when the number of centroids becomes large and KHATRI-RAO- k -MEANS can become slower than the naïve approach to Khatri-Rao k -Means. On the other hand, standard k -MEANS is consistently the fastest algorithm.

As for memory requirements, KHATRI-RAO- k -MEANS often incurs similar memory requirements as standard k -MEANS that uses $h_1 + h_2$ vectors to represent centroids. On the other hand, k -MEANS with $h_1 h_2$ centroids can be more memory demanding, particularly when the number of centroids grows. Specifically, as the number of centroids increases, k -MEANS with $h_1 h_2$ centroids uses up to about 2.7 times as much memory as KHATRI-RAO- k -MEANS in both the BLOBS and CLASSIFICATION datasets. Furthermore, the gap grows with the number of centroids, corroborating the space-complexity discussion in Section 6.

It is also of interest to investigate the scalability of KHATRI-RAO- k -MEANS as a function of the number of protocentroids sets. For the same experimental settings as in Figure 7, Figure 12 reports

the runtime and memory usage of KHATRI-RAO- k -MEANS as the number p of protocentroids sets varies in {2, 3, 4}, while keeping a budget of 12 vectors to represent the centroids. As suggested by the results, increasing the number of protocentroids leads to higher runtime, as more centroids are computed. Memory usage may also increase slightly, with variations always bounded by 5 mebibytes. For reference, we additionally report the runtime and peak memory usage of all baseline methods with $h_1 = h_2$. These measurements are obtained on the original BLOBS and CLASSIFICATION datasets, corresponding to relatively small problem instances.

Impact of the number of centroids. To complement the preliminary experiments presented in Section 9, Figure 13 shows the inertia incurred by the KHATRI-RAO- k -MEANS and the baselines under comparison as the number of ground-truth clusters k in the data is varied within {25, 100, 225, 400, 625} in the BLOBS and CLASSIFICATION datasets. The figure demonstrates that KHATRI-RAO- k -MEANS consistently achieves significantly lower inertia than the baselines at parity parameters (i.e., number of vectors to represent centroids). Standard k -MEANS with $h_1 h_2$ vectors to represent centroids achieves the lowest inertia, although its advantage over KHATRI-RAO- k -MEANS, which uses only $h_1 + h_2$ vectors, remains consistently small.

Finally, the naïve approach always incurs the highest inertia, and its performance tends to worsen with the number of clusters, which is not the case for the other algorithms under comparison.

Impact of the reparameterization on the autoencoder. As explained in Section 4.3, in Khatri-Rao deep clustering, we compress the autoencoder by reparameterizing its weights as Hadamard product of low-rank factors.

This reparameterization can impact training stability and clustering performance. Thus, in the following, we empirically investigate the training stability and clustering performance as a function of the rank of both Hadamard factors. In particular, the rank of the Hadamard factors is varied in {20, 30, 40, 50, 60}.

Figure 14 displays the value of the clustering quality metrics monitored in this work for the BLOBS and CLASSIFICATION datasets against the rank of each of two Hadamard factors. As the figure suggests, the chosen rank can have an impact on the clustering performance of Khatri-Rao deep clustering. However, such impact is limited. Increasing the rank of both Hadamard factors from 10 to 60 results in changes of at most 0.1 across all metrics, with most variations remaining below 0.05. On the BLOBS dataset, the smallest rank considered (20) yields better clustering performance than higher ranks. In this case, the implicit regularization induced by the reparameterization as Hadamard product of low-rank factors may be beneficial for clustering.

For the BLOBS and CLASSIFICATION datasets, Figure 15 depicts the evolution of the training loss over 500 epochs for both the uncompressed autoencoder and the compressed autoencoder in case the rank of both Hadamard factors is 20, 40 and 60. Although, particularly for the lower ranks, the compressed autoencoder exhibits a slightly slower initial decrease in training loss compared to the uncompressed autoencoder, it quickly converges to a comparable loss level. Overall, the training loss decays smoothly and consistently over iterations for both the uncompressed and compressed autoencoders across all considered ranks.

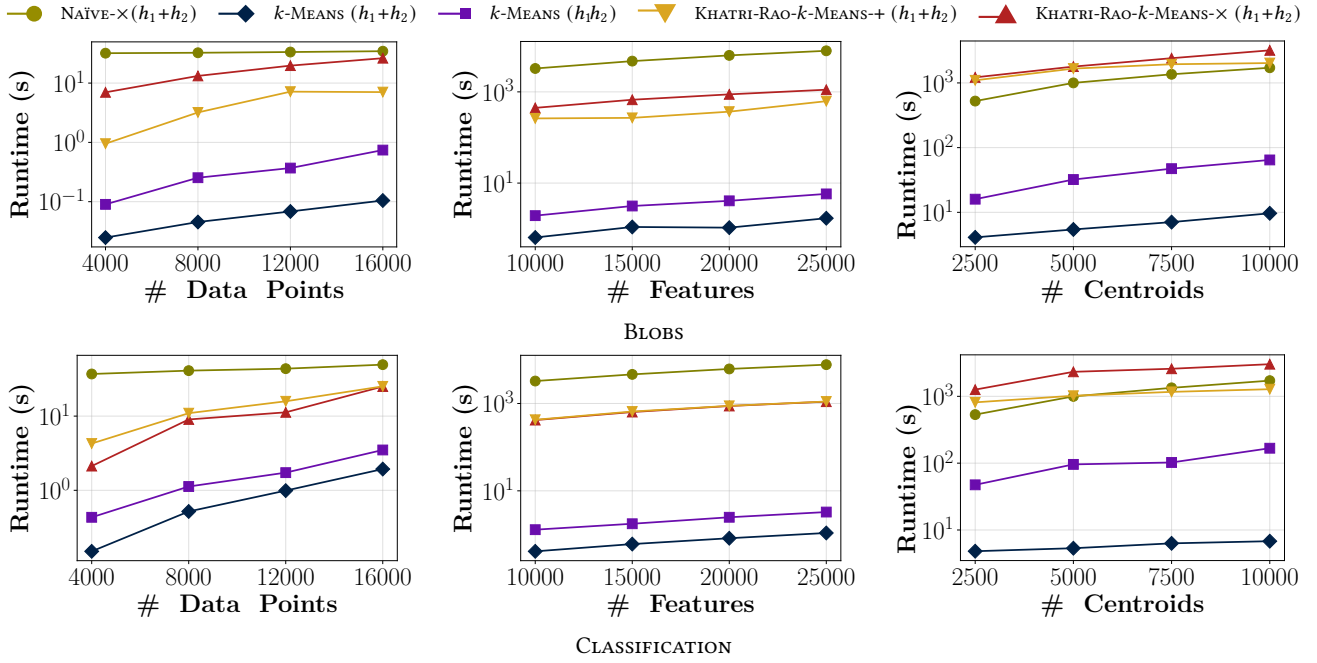


Figure 10: Experiments on BLOBS and CLASSIFICATION datasets. Runtime (in seconds) by number of data points n , number of features m and number of centroids k . Here, k -MEANS (h_1h_2) uses h_1h_2 vectors to represent centroids, while all other algorithms h_1+h_2 . The y-axis is on a logarithmic scale.

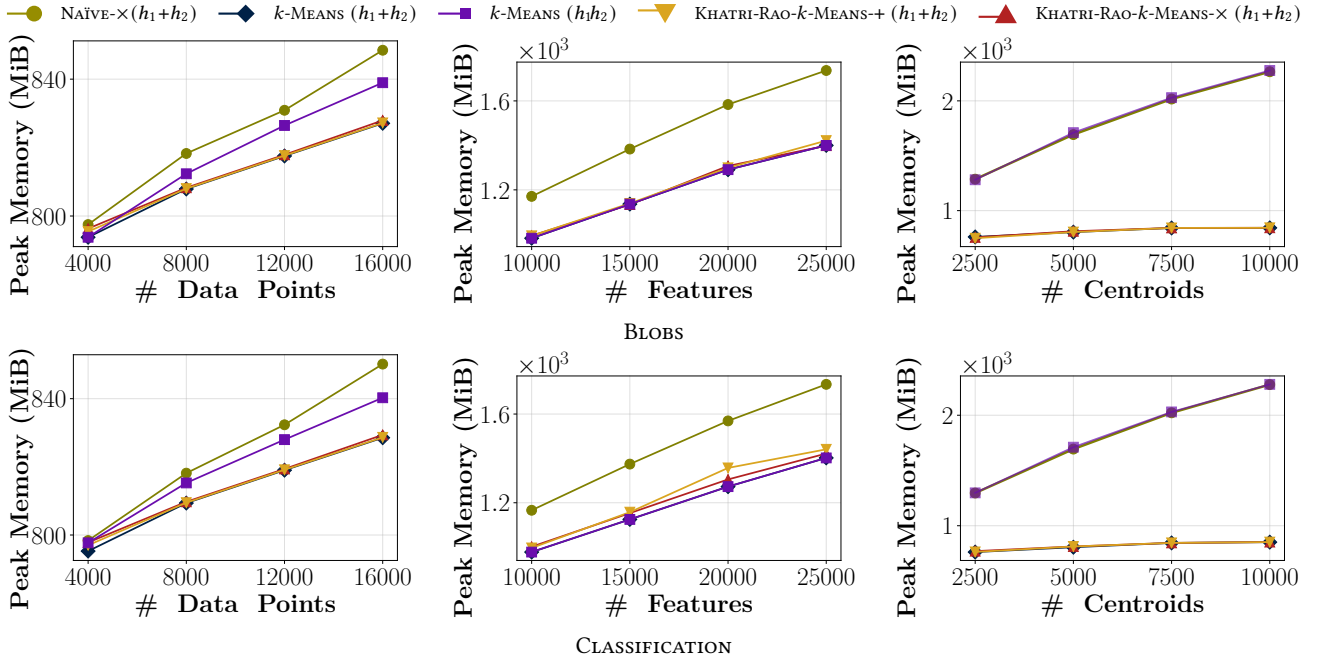


Figure 11: Experiments on BLOBS and CLASSIFICATION datasets. Peak memory usage (in Megabytes) by number of data points n , number of features m and number of centroids k . Here, k -MEANS (h_1h_2) uses h_1h_2 vectors to represent centroids, while all other algorithms h_1+h_2 . The y-axis is on a logarithmic scale.

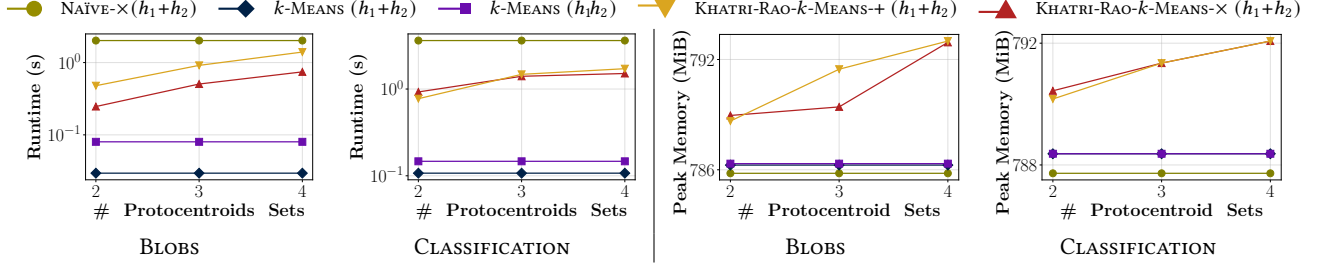


Figure 12: Experiments on BLOBS and CLASSIFICATION datasets. Runtime in seconds (left) and memory usage in Megabytes (right) by number of sets of protocentroids. Here, KHATRI-RAO-k-MEANS always uses 12 total vectors to represent centroids. All baselines assume $h_1 = h_2 = 6$. The y-axis is on a logarithmic scale.

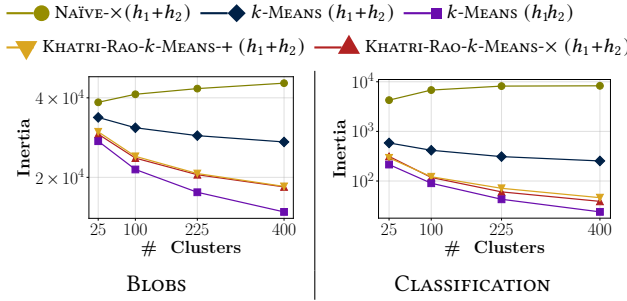


Figure 13: Experiments on BLOBS and CLASSIFICATION datasets. Inertia as a function of the number of ground-truth clusters k . Here, k -MEANS (h_1h_2) uses h_1h_2 vectors to represent centroids, while all other algorithms $h_1 + h_2$. The y-axis is on a logarithmic scale.

Appendix D DETAILS ON KHATRI-RAO CLUSTERING AND MATRIX DECOMPOSITION

As anticipated in Section 2, clustering is closely related to matrix decomposition. Both centroid-based clustering and matrix decomposition pursue the objective of summarizing data through a succinct set of prototypes. Many clustering problems such as k -Means can be seen as a particular case of matrix decomposition where one of the factor matrices is constrained to indicate cluster membership. In fact, it has been shown that the continuous relaxation of the k -Means cluster-membership indicators is given by the principal components of the data [7], establishing a direct link between k -Means clustering and principal component analysis, a prominent matrix-decomposition technique. Clustering via matrix decomposition has been successfully applied in different domains like text mining [22], graph mining [1, 13] and bioinformatics [2].

The established relationship between clustering and matrix factorization also extends to tensor data. Tensor decompositions can be used for clustering tensor data [14], with successful applications in areas like text mining across multiple languages [4], multilayer graph mining [3] and bioinformatics [21].

Given the close interplay between matrix (or tensor) decomposition and clustering, a natural question is how Khatri-Rao clustering fits within the broader framework of matrix decomposition.

In this section, we first review how standard k -Means clustering can be formulated as a constrained matrix-decomposition problem. Second, we show that also Khatri-Rao k -Means can be interpreted as a constrained instance of matrix decomposition. Then, we discuss the well-known example of nonnegative matrix factorization [18]. Finally, we present a simple illustrative experiment which anticipates the potential advantages of Khatri-Rao clustering in clustering via matrix decomposition.

k -Means clustering as matrix decomposition. The k -Means clustering problem can be formulated as approximating a data matrix by a product of lower-rank factors: one factor encodes cluster centroids and the other encodes cluster-membership assignments.

More specifically, the k -Means problem admits a natural reformulation as a constrained matrix decomposition task as follows. Let $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^m$ be the dataset and denote by $\Theta = \{\mu_i\}_{i=1}^k$ the set of centroids, where μ_i is the centroid of C_i . Introduce the assignment matrix $\mathbf{Z} \in \{0, 1\}^{n \times k}$, where $Z_{ji} = 1$ if $\mathbf{x}_j \in C_i$ and 0 otherwise, so that each data point is assigned to exactly one cluster (i.e., $\mathbf{Z}_j \mathbf{1}_k = \mathbf{1}_n$). Let $\mathbf{C} \in \mathbb{R}^{k \times m}$ be the centroid matrix whose i -th row equals μ_i^\top . Collecting the data points in \mathcal{D} as the rows of a data matrix $\mathbf{D} \in \mathbb{R}^{n \times m}$, the inertia objective in Equation (1), minimized in k -Means, can be written compactly as:

$$Q_C(\mathcal{D}, \Theta) = \|\mathbf{D} - \mathbf{Z}\mathbf{C}\|_F^2 \quad (9)$$

under the aforementioned assignment constraints encoded in \mathbf{Z} . Here, $\|\mathbf{X}\|_F$ denotes the Frobenius norm of \mathbf{X} . Thus, the k -Means problem is equivalent to seeking a decomposition $\mathbf{D} \approx \mathbf{Z}\mathbf{C}$ under the constraint that \mathbf{Z} is a binary assignment matrix with one nonzero entry per row.

Khatri-Rao k -Means clustering as constrained matrix decomposition. Like standard k -Means clustering, Khatri-Rao k -Means clustering also admits a formulation in terms of matrix decomposition. To derive the formulation of Khatri-Rao k -Means as matrix decomposition, let us introduce matrices \mathbf{Z}_i and \mathbf{P}_i with $i = 1 \dots p$. Similarly to \mathbf{Z} for standard k -Means, matrices \mathbf{Z}_i have one-hot assignment vectors as rows. Matrix \mathbf{Z}_i assigns each data point to one of h_i protocentroids stored as rows of matrix \mathbf{P}_i . Then, the Khatri-Rao k -Means problem amounts to the minimization of:

$$Q_C(\mathcal{D}, \Theta) = \|\mathbf{D} - (\mathbf{Z}_1\mathbf{P}_1) \oplus (\mathbf{Z}_2\mathbf{P}_2) \dots \oplus (\mathbf{Z}_p\mathbf{P}_p)\|_F^2. \quad (10)$$

under the assignment constraints in matrices \mathbf{Z}_i . Hence, Khatri-Rao k -Means clustering induces a decomposition of the form $\mathbf{D} \approx$

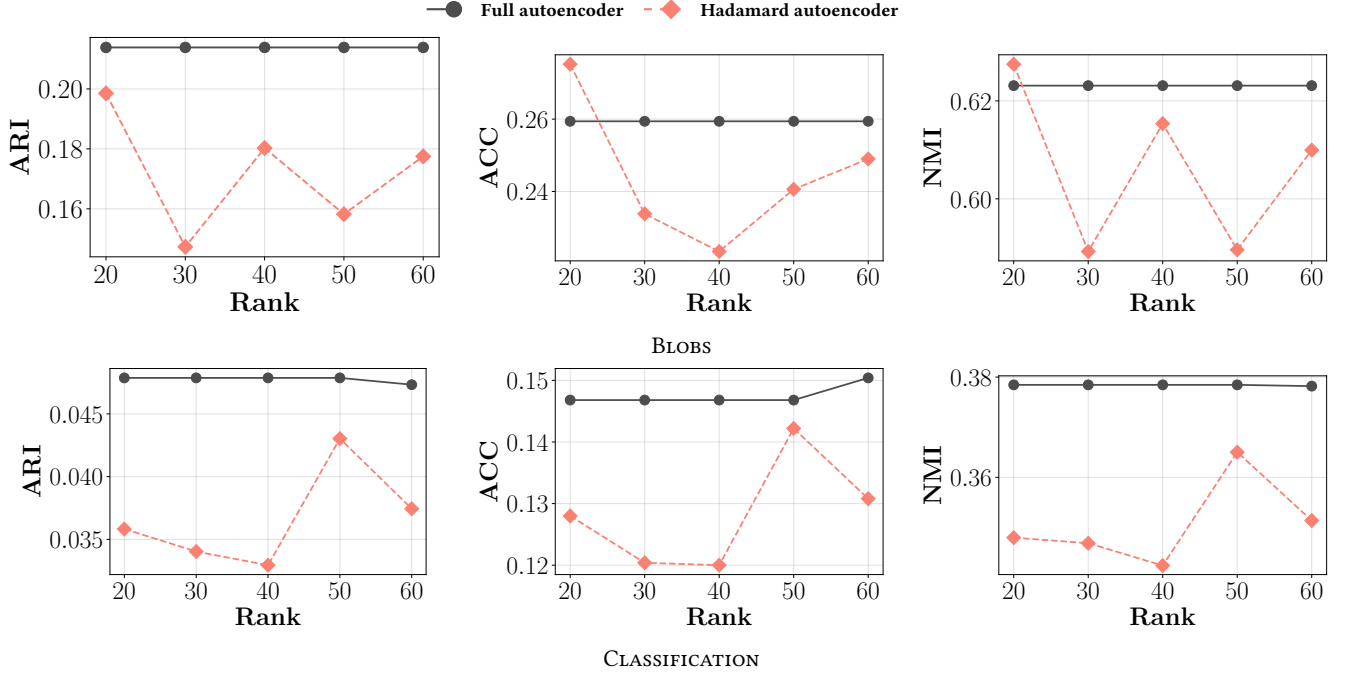


Figure 14: Experiments on BLOBS and CLASSIFICATION datasets. Unsupervised clustering accuracy (ACC), adjusted Rand index (ARI) and normalized mutual information (NMI) for KHATRI-RAO DKM against rank of both factors in the reparameterization of the autoencoder weights as Hadamard product of two low-rank factors.

$(Z_1 P_1) \oplus (Z_2 P_2) \dots \oplus (Z_p P_p)$ under the constraint that Z_i are binary assignment matrices with one nonzero entry per row. Such decomposition is illustrated in Figure 16.

Khatri-Rao clustering and the Hadamard decomposition. Equation (10) highlights the connection between Khatri-Rao k -Means and the Hadamard decomposition, studied in our previous work [6]. In the Hadamard decomposition, as in general matrix decomposition, the factors Z_i are unconstrained. As a result, summation does not introduce additional structure: the sum of q rank- r factors can represent matrices of rank at most qr , and summing two decompositions each using h vectors per factor is, in principle, equivalent to a single decomposition with $2h$ vectors per factor. In the setting of matrix decomposition, summation of multiple matrix decompositions merely reallocates expressiveness. By contrast, the Hadamard product can yield a more expressive model since the product of q rank- r factors can represent matrices of rank up to r^q , which explains the improved compression rates observed for the Hadamard decomposition in previous work [6, 10]. This expressiveness advantage is also the reason why we adopt the Hadamard product in the reparameterization of the autoencoder for Khatri-Rao deep clustering.

The expressiveness advantage of the Hadamard product over the sum in matrix decomposition, however, does not carry over to Khatri-Rao clustering. Once the factors Z_i are constrained to be one-hot cluster indicators, multiple decompositions combined additively are no longer equivalent in principle to a single enlarged factor. The discrete assignment structure prevents expressiveness

from being freely reallocated, so that summation can introduce genuinely new expressiveness.

The example of nonnegative matrix factorization. Nonnegative matrix factorization (NMF) exemplifies how clustering and matrix decomposition are fundamentally intertwined.

Given an integer k , NMF represents data as a product of nonnegative factors [16], imposing the model $D = WH$ with $W_{i,j} \geq 0 \forall i = 1 \dots n, j = 1 \dots k$ and $H_{i,j} \geq 0 \forall i = 1 \dots k, j = 1 \dots m$. NMF has emerged as a powerful approach to clustering, as its representations naturally induce groupings of data points and features.

It is known that optimizing NMF under the orthogonality constraint $W^T W = I$, i.e., solving:

$$\min_{W, H; W_{i,j} \geq 0 \forall i, j, H_{i,j} \geq 0 \forall i, j, W^T W = I} \|D - WH\|_F^2 \quad (11)$$

is equivalent to the k -MEANS problem (assuming positive data) [18].

As in the case of standard clustering, a Khatri-Rao variant of NMF asks to solve:

$$\min_{\substack{W_t, H_t \\ W_{t,i,j} \geq 0, H_{t,i,j} \geq 0, \forall t, i, j \\ W_t^T W_t = I \forall t}} \|D - (W_1 H_1) \oplus (W_2 H_2) \dots \oplus (W_p H_p)\|_F^2 \quad (12)$$

Algorithms that solve the optimization problem in Equation (11) have been proposed in the literature [5]. Studying their extension to

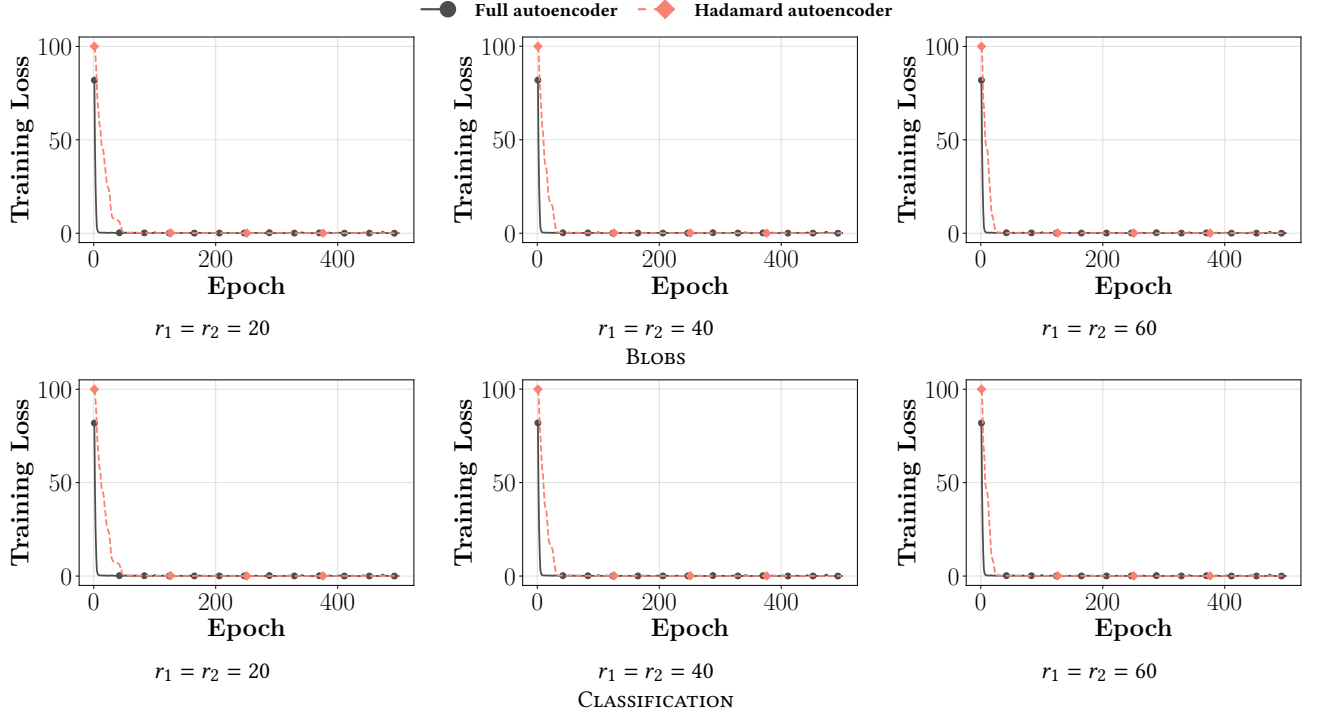


Figure 15: Experiments on BLOBS and CLASSIFICATION datasets. Training loss decay across 500 epochs for the unconstrained autoencoder and the autoencoder with weights reparameterized as the Hadamard product of two low-rank factors with ranks r_1 and r_2 . We set $r_1 = r_2$, and we show results for three different choices of r_1 and r_2 (20, 40 and 60).

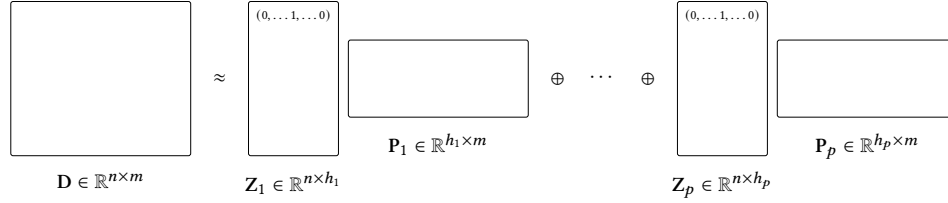


Figure 16: Diagram representing Khatri-Rao k -Means as matrix decomposition.

the Khatri-Rao clustering paradigm by addressing the optimization problem in Equation (12) is an exciting direction of future work.

D.1 Illustrative experiments.

To conclude, we carry out an illustrative experiment to anticipate empirically the benefits that Khatri-Rao clustering can provide to centroid-based clustering via matrix decomposition. This experiment is not intended to introduce a new state-of-the-art clustering method based on matrix decomposition, but to demonstrate the practical advantages that the Khatri-Rao clustering paradigm can potentially offer in this context through a simple approach. We focus on the case of two sets of prototypical centroids. Equations (9) and (10) express the k -Means and Khatri-Rao k -Means problems in terms of matrix decomposition, respectively. In this illustrative experiment, we address such problems via gradient-based optimization, relying on PyTorch for automatic differentiation. In particular, for

k -Means clustering, we define the model $D = ZC$, and we learn its parameters by optimizing the following loss function:

$$\sum_{i=1}^n \sum_{j=1}^k z_{i,j} \|D_i - C_j\|^2. \quad (13)$$

Similarly, for Khatri-Rao clustering, we define the model $D = (Z_1 P_1) \oplus (Z_2 P_2)$, and we learn its parameters by optimizing the loss function in Equation (13) adapted to enforce the Khatri-Rao structure. For optimization purposes, we relax the constraints that matrices Z in standard clustering and Z_i in Khatri-Rao clustering are binary assignment matrices. Instead, each row encodes continuous soft assignments to all clusters. We run the optimization for 5000 epochs using the Adam optimizer, with gradients computed over the entire dataset at each epoch. To obtain one-hot assignment vectors, we first initialize the continuous soft assignments at random, and we apply the softmax function to ensure that assignments for each

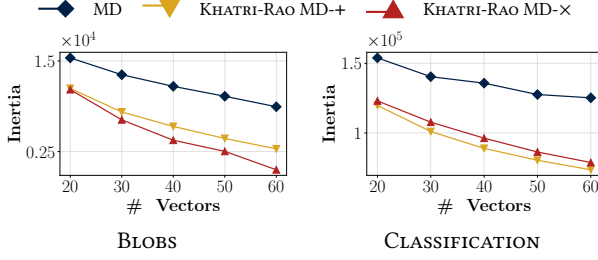


Figure 17: Experiments on BLOBS and CLASSIFICATION datasets. Inertia by number of vectors used to represent centroids for a simple approach to standard clustering via matrix decomposition (MD) and its Khatri-Rao variant with sum aggregator (KHATRI-RAO MD+) and product aggregator (KHATRI-RAO MD-x). The y-axis is on a logarithmic scale.

data point sum to 1. During the optimization, the temperature of the softmax function is linearly annealed from 1 to 0.1 to progressively sharpen the assignments until they approach one-hot vectors. The learning rate is also annealed at the same time from 10^{-1} to 5×10^{-3} . More specifically, both temperature and learning rate follow linear schedules of the form $(1 - a_t)l_{start} + a_t l_{end}$, where l_{start} and l_{end} are either the initial and final temperature or learning rate and a_t increases linearly over epochs. To improve performance, we initialize the centroids (stored as rows of matrix C for standard clustering and as combinations of rows of matrices P_i for Khatri-Rao clustering) uniformly in the range of the data.

Upon termination of the 5000 epochs, we compute the hard inertia of the solution clustering by assigning each data point to the cluster with the highest assignment score. Because of the temperature annealing, highest assignment scores are typically close to 1 while the others approach 0.

The results of the described experiment for the BLOBS and CLASSIFICATION datasets are displayed in Figure 17. Here, we compare clustering via standard matrix decomposition (MD) and its Khatri-Rao clustering extension with sum aggregator (KHATRI-RAO MD+) and product aggregator (KHATRI-RAO MD-x), as the number of vectors used to represent centroids varies in $\{20, 30, 40, 50, 60\}$. The results of the experiment indicate that the matrix decomposition approach to standard clustering incurs higher inertia than its Khatri-Rao counterpart. More specifically, standard matrix decomposition exhibits an inertia that is at least 83% and 28% higher than that achieved by any Khatri-Rao variant in the BLOBS and CLASSIFICATION datasets, respectively. The advantages provided by the Khatri-Rao clustering paradigm also tend to increase with the number of vectors used to represent centroids.

Appendix E PROOFS

In this section, we collect the proofs of the propositions stated in the paper.

Proof of Proposition 6.1.

PROOF. We illustrate the proof for the product aggregator. Consider the j -th protocentroid in the first set of protocentroids. The

optimal update for this protocentroid satisfies:

$$\theta_1^j = \arg \min_{\theta} \sum_{l=1}^{h_2} \sum_{\mathbf{x} \in C_{j,l}} (\mathbf{x} - \theta \odot \theta_2^l)^2.$$

Therefore θ_1^j can be found by computing the gradient of this sum of squared differences with respect to θ and equating it to the zero vector $\mathbf{0}$. Doing so, one obtains:

$$-2 \sum_{l=1}^{h_2} \sum_{\mathbf{x} \in C_{j,l}} (\mathbf{x} - \theta \odot \theta_2^l) \theta_2^l = \mathbf{0},$$

which holds if and only if:

$$\theta_1^j = \frac{\sum_{l=1}^{h_2} \sum_{\mathbf{x} \in C_{j,l}} \mathbf{x} \odot \theta_2^l}{\sum_{l=1}^{h_2} |C_{j,l}| \theta_2^l \odot \theta_2^l}.$$

The derivations for the second set of protocentroids are symmetric, and the proof for the sum aggregator is also similar. \square

Proof of Proposition 8.1.

PROOF. Assuming p protocentroid sets of equal size and exact budget usage, each set contains $h = \frac{b}{p}$ protocentroids, so the total number of centroids that can be represented is:

$$\left(\frac{b}{p}\right)^p.$$

Maximizing this function over $p > 0$ is equivalent to maximizing its logarithm:

$$p \log\left(\frac{b}{p}\right) = p \log b - p \log p,$$

where \log denotes the natural logarithm. Differentiating with respect to p yields:

$$\log\left(\frac{b}{p}\right) - 1,$$

which is positive for $p < \frac{b}{e}$ and negative for $p > \frac{b}{e}$. Moreover, the second derivative:

$$-\frac{1}{p} < 0 \quad (p > 0)$$

shows that the log-objective is strictly concave, and therefore attains a unique maximum at $p = \frac{b}{e}$ in the continuous domain.

When restricting to the divisor-only setting, the optimal value of p , p^{max} , is required to be an integer that exactly divides b , so that $h^{max} = \frac{b}{p^{max}}$ is an integer. Since the objective is strictly increasing for $p < \frac{b}{e}$ and strictly decreasing for $p > \frac{b}{e}$, among all such admissible values the maximum is attained at one of the two values of p that are immediately below or above $\frac{b}{e}$ (when each exists). All other admissible divisors are necessarily further away from $\frac{b}{e}$, and therefore yield a strictly smaller objective value. \square

Proof of Proposition 8.2.

PROOF. We present the proof for the product aggregator. However, extension to different aggregator functions is straightforward.

First, to see why it must be $p^* \geq \log_{h_{min}} k$, notice that, if $p^* < \log_{h_{min}} k$ then $h_{min}^{p^*} < k$, and $h_{min}^{p^*}$ is the maximum number of centroids that can be represented using p^* sets of at least h_{min}

protocentroids. As regards the other inequality, to prove that $p^* \leq \lceil \frac{k}{h_{\min}-1} \rceil$, it is sufficient to construct an example where all sets have a protocentroid with all entries equal to 1 and the remaining at least $h_{\min} - 1$ protocentroids that are equal to centroids. Different sets contain different centroids. In the illustrated scenario, each set of protocentroids represents exactly at least $h_{\min} - 1$ centroids. Thus, $\lceil \frac{k}{h_{\min}-1} \rceil$ sets of protocentroids can always represent k centroids. \square

REFERENCES

- [1] Kamal Berahmand, Mehrnosh Mohammadi, Razieh Sheikhpour, Yuefeng Li, and Yue Xu. 2024. WSNMF: Weighted Symmetric Nonnegative Matrix Factorization for attributed graph clustering. *Neurocomputing* 566 (2024), 127041.
- [2] Jean-Philippe Brunet, Pablo Tamayo, Todd R Golub, and Jill P Mesirov. 2004. Metagenes and molecular pattern discovery using matrix factorization. *Proceedings of the national academy of sciences* 101, 12 (2004), 4164–4169.
- [3] Zitai Chen, Chuan Chen, Zibin Zheng, and Yi Zhu. 2019. Tensor Decomposition for Multilayer Networks Clustering. In *AAAI*. AAAI Press, 3371–3378.
- [4] Peter A. Chew, Brett W. Bader, Tamara G. Kolda, and Ahmed Abdelali. 2007. Cross-language information retrieval using PARAFAC2. In *KDD*. ACM, 143–152.
- [5] Seungjin Choi. 2008. Algorithms for orthogonal nonnegative matrix factorization. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*. IEEE, 1828–1832.
- [6] Martino Ciaperoni, Aristides Gionis, and Heikki Mannila. 2024. The Hadamard decomposition problem. *Data Mining and Knowledge Discovery* (2024), 1–42.
- [7] Chris Ding and Xiaofeng He. 2004. K-means clustering via principal component analysis. In *Proceedings of the twenty-first international conference on Machine learning*. 29.
- [8] Marek Gagolewski. 2022. A framework for benchmarking clustering algorithms. *SoftwareX* 20 (2022), 101270.
- [9] Stephan Günnemann, Ines Färber, Matthias Sebastian Rüdiger, and Thomas Seidl. 2014. SMVC: semi-supervised multi-view clustering in subspace projections. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014*. ACM, 253–262. <https://doi.org/10.1145/2623330.2623734>
- [10] Nam Hyeon-Woo, Moon Ye-Bin, and Tae-Hyun Oh. 2021. FedPara: Low-rank Hadamard Product for Communication-Efficient Federated Learning. In *International Conference on Learning Representations*.
- [11] George Karypis, Eui-Hong Han, and Vipin Kumar. 1999. Chameleon: Hierarchical Clustering Using Dynamic Modeling. *Computer* 32, 8 (1999), 68–75. <https://doi.org/10.1109/2.781637>
- [12] Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. 2023. The UCI machine learning repository.
- [13] Mikhail Khodak, Neil A. Tenenholz, Lester Mackey, and Nicolò Fusi. 2021. Initialization and Regularization of Factorized Neural Layers. In *ICLR*. OpenReview.net.
- [14] Tamara G. Kolda and Brett W. Bader. 2009. Tensor Decompositions and Applications. *SIAM Rev.* 51, 3 (2009), 455–500.
- [15] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [16] Daniel Lee and H Sebastian Seung. 2000. Algorithms for non-negative matrix factorization. *Advances in neural information processing systems* 13 (2000).
- [17] Collin Leiber, Lukas Miklautz, Claudia Plant, and Christian Böhm. 2023. Benchmarking deep clustering algorithms with clustpy. In *2023 IEEE International Conference on Data Mining Workshops (ICDMW)*. IEEE, 625–632.
- [18] Tao Li and Cha-charis Ding. 2018. Nonnegative matrix factorizations for clustering: A survey. *Data Clustering* (2018), 149–176.
- [19] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch. (2017).
- [20] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [21] Matteo Ruffini, Ricard Gavaldà, and Esther Limon. 2017. Clustering Patients with Tensor Decomposition. In *MLHC (Proceedings of Machine Learning Research, Vol. 68)*. PMLR, 126–146.
- [22] Wei Xu, Xin Liu, and Yihong Gong. 2003. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*. 267–273.