



Scalable Unsupervised Feature Selection with Reconstruction Error Guarantees via QMR Decomposition

Ciwan Ceylan
ciwan@kth.se

KTH Royal Institute of Technology
Stockholm, Sweden

Kambiz Ghoorchian
kambiz.ghoorchian@seb.se

SEB Group
Stockholm, Sweden

Danica Kragic
dani@kth.se

KTH Royal Institute of Technology
Stockholm, Sweden

Abstract

Unsupervised feature selection (UFS) methods have garnered significant attention for their capability to eliminate redundant features without relying on class label information. However, their scalability to large datasets remains a challenge, rendering common UFS methods impractical for such applications. To address this issue, we introduce QMR-FS, a greedy forward filtering approach that selects linearly independent features up to a specified relative tolerance, ensuring that any excluded features can be reconstructed from the retained set within this tolerance. This is achieved through the QMR matrix decomposition, which builds upon the well-known QR decomposition. QMR-FS benefits from linear complexity relative to the number of instances and boasts exceptional performance due to its ability to leverage parallelized computation on both CPU and GPU. Despite its greedy nature, QMR-FS achieves comparable classification and clustering accuracies across multiple datasets when compared to other UFS methods, while achieving runtimes approximately 10 times faster than recently proposed scalable UFS methods for datasets ranging from 100 million to 1 billion elements.

CCS Concepts

• **Computing methodologies** → **Feature selection; Unsupervised learning.**

Keywords

Unsupervised learning, Feature selection, Scalability, Linear independence

ACM Reference Format:

Ciwan Ceylan, Kambiz Ghoorchian, and Danica Kragic. 2024. Scalable Unsupervised Feature Selection with Reconstruction Error Guarantees via QMR Decomposition. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management (CIKM '24)*, October 21–25, 2024, Boise, ID, USA. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3627673.3679994>

1 Introduction

Unsupervised feature selection (UFS) involves extracting a subset of features from a dataset by eliminating redundant ones without relying on task-specific information, such as class labels. The removal of redundant features offers several inherent benefits: it facilitates

data visualization and understanding, reduces measurement and storage requirements, and decreases training and utilization times [10, 25]. Additionally, feature selection can enhance the performance of downstream tasks, such as clustering, by mitigating the curse of dimensionality [10, 13, 22].

UFS methods are categorized into two main types: *wrapper* and *filter* methods [25]. Wrapper methods rely on an external task algorithm, such as a clustering algorithm, to evaluate feature importance, whereas filter methods use only the intrinsic properties of the features and dataset. Filter methods are considered faster and more scalable [10, 25, 26], making them the focus of this work.

Despite being labelled 'fast', filtering UFS methods often become inefficient for large datasets [15] due to their quadratic time complexity with respect to the number of instances [11, 18, 31–34], meaning n in an $n \times d$ data matrix with d features. This quadratic complexity makes computation intractable for the datasets that would benefit most from feature selection in terms of reducing storage and utilization times. While some methods exhibit linear complexity in n [29, 30], they use iterative schemes that must run until convergence, which is slow in practice.

To address the lack of scalable UFS methods, we introduce QMR-FS, a fast feature selection algorithm. QMR-FS employs greedy forward feature selection, where features are considered one-by-one in a fixed order, and feature redundancy is measured only with respect to already selected features. Specifically, QMR-FS uses linear reconstruction error to measure redundancy, based on the premise that linearly reconstructable features are redundant [6]. Our QMR decomposition, an extension of the QR decomposition [12, Ch. 2.1], computes the reconstruction efficiently by avoiding costly linear model fitting, ensuring linear complexity in n , and is optimized for parallel implementation on both CPU and GPU. Additionally, QMR-FS guarantees that reconstruction errors are smaller than a specified tolerance relative to the feature columns' L2-norm, allowing it to retain all features if redundancy criteria are not met. Thus, QMR-FS is suitable for scenarios where the potential for feature removal without information loss is uncertain.

We demonstrate the scalability of QMR-FS via runtimes around 30 seconds on datasets with up to 7.88 million instances and a total of 1 billion elements, roughly 10 times faster than recent scalable UFS methods [15]. Furthermore, the features selected by QMR-FS exhibit comparable classification accuracy and clustering normalized mutual information (NMI) scores relative to both common and recent UFS methods, which require significantly longer runtimes. Code for QMR-FS and our experiments is available online¹.



This work is licensed under a Creative Commons Attribution International 4.0 License.

CIKM '24, October 21–25, 2024, Boise, ID, USA
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0436-9/24/10
<https://doi.org/10.1145/3627673.3679994>

¹<https://github.com/ciwanceylan/qmr-feature-selection>

This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation.

2 Method

We now describe how our QMR-FS method performs scalable feature selection on a data matrix with n instances and d features, $\mathbf{X} \in \mathbb{R}^{n \times d}$. Since our focus is on scalability to large n , we limit this paper to the case $n \geq d$ and relegate the opposite case to future work.

QMR-FS performs greedy forward feature selection, considering features from left to right in \mathbf{X} , with the decision to retain a feature based on the previously retained features. The guiding principle is that of linear independence. We say that a feature column $\mathbf{x}_j \in \mathbf{X}$ is *linearly dependent* on previous feature columns if it can be expressed as a linear combination of them, plus a constant term,

$$\mathbf{x}_j = \sum_{k=1}^{j-1} a_k \mathbf{x}_k + b \mathbf{1}_n, \quad (1)$$

where a_k and b are coefficients, and $\mathbf{1}_n$ is a length n vector of all ones. This is expressed concisely as $\mathbf{x}_j \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{1}_n)$. Conversely, if such a_k and b do not exist, \mathbf{x}_j is *linearly independent*, $\mathbf{x}_j \notin \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{1}_n)$, and should be retained.

To determine the independence of \mathbf{x}_j without having to explicitly find a_k and b through expensive linear model fitting, we make use of the following lemma, central to our method.

LEMMA 2.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{R}_{\text{ref}} \in \mathbb{R}^{d \times d}$ satisfy $\mathbf{X} = \mathbf{U} \mathbf{R}_{\text{ref}}$ and $\mathbf{R}_{\text{ref}} = \mathbf{U}^\dagger \mathbf{X}$, where \mathbf{R}_{ref} is in row echelon form (REF), and $\mathbf{U} \in \mathbb{R}^{n \times d}$ has full rank and the left-side inverse $\mathbf{U}^\dagger \in \mathbb{R}^{d \times n}$. Then $\mathbf{x}_j \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$ iff $\mathbf{r}_j \in \text{span}(\mathbf{r}_1, \dots, \mathbf{r}_{j-1})$ for any $j > 1$ where $\mathbf{x}_j \in \mathbb{R}^n$ is the j th column in \mathbf{X} , and $\mathbf{r}_j \in \mathbb{R}^d$ is the j th column in \mathbf{R}_{ref} .*

PROOF. We first show that $\mathbf{x}_j \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1}) \implies \mathbf{r}_j \in \text{span}(\mathbf{r}_1, \dots, \mathbf{r}_{j-1})$. In this case, we can assume $\exists a_{i,j} \in \mathbb{R}$ s.t. $\mathbf{x}_j = \sum_{i=1}^{j-1} a_{i,j} \mathbf{x}_i$. From $\mathbf{X} = \mathbf{U} \mathbf{R}_{\text{ref}}$, we know that $\mathbf{x}_i = \mathbf{U} \mathbf{r}_i$. Thus,

$$\mathbf{U} \mathbf{r}_j = \mathbf{x}_j = \sum_{i=1}^{j-1} a_{i,j} \mathbf{x}_i = \sum_{i=1}^{j-1} a_{i,j} \mathbf{U} \mathbf{r}_i = \mathbf{U} \sum_{i=1}^{j-1} a_{i,j} \mathbf{r}_i.$$

Multiplication with \mathbf{U}^\dagger completes the first part of the proof,

$$\mathbf{r}_j = \mathbf{U}^\dagger \mathbf{U} \sum_{i=1}^{j-1} a_{i,j} \mathbf{r}_i = \sum_{i=1}^{j-1} a_{i,j} \mathbf{r}_i.$$

Next, $\mathbf{r}_j \in \text{span}(\mathbf{r}_1, \dots, \mathbf{r}_{j-1}) \implies \mathbf{x}_j \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$ is proved with a block matrix representation of $\mathbf{X} = \mathbf{U} \mathbf{R}_{\text{ref}}$ as aid,

$$\left(\begin{array}{c|c|c|c|c} \mathbf{x}_1 & \dots & \mathbf{x}_{j-1} & \mathbf{x}_j & \dots & \mathbf{x}_d \\ \hline \mathbf{u}_1 & \dots & \mathbf{u}_h & \dots & \mathbf{u}_d & \end{array} \right) = \left(\begin{array}{c|c|c|c|c} \mathbf{u}_1 & \dots & \mathbf{u}_h & \dots & \mathbf{u}_d & \end{array} \right) \begin{pmatrix} r_{1,1} & r_{1,2} & \dots & r_{1,j-1} & r_{1,j} & \dots & r_{1,d} \\ r_{2,2} & \dots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ r_{h,j-1} & r_{h,j} & \dots & r_{h,j} & \dots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix}. \quad (2)$$

Here, h is the rank of the first $j-1$ columns in \mathbf{X} , which coincide with the number of rows with pivot elements before column j in \mathbf{R}_{ref} . Since \mathbf{R}_{ref} is in REF, $h \leq j-1$, and by the assumption $\mathbf{r}_j \in \text{span}(\mathbf{r}_1, \dots, \mathbf{r}_{j-1})$ we know that \mathbf{r}_j cannot have pivot elements.

We know that $\mathbf{x}_j = \sum_{i=1}^h \mathbf{u}_i r_{i,j}$ (orange and brown in Eq. (2)), and we wish to express \mathbf{u}_i using columns in \mathbf{X} . To do so, we look at the block equation highlighted in blue and brown blocks in Eq. (2),

$$\mathbf{X}_{:,1:j-1} \mathbf{P} = \mathbf{U}_{:,1:h} \mathbf{R}_{\text{ref}1:h,1:j-1}. \quad (3)$$

If $h > 0$, the block matrix $\mathbf{R}_{\text{ref}1:h,1:j-1} \in \mathbb{R}^{h \times j-1}$ has pivot elements in every row by construction (REF), meaning that it has

full row rank. Therefore, it has the right-side inverse $\mathbf{P} \in \mathbb{R}^{j-1 \times h}$ such that $\mathbf{R}_{\text{ref}1:h,1:j-1} \mathbf{P} = \mathbf{I}_{h \times h}$ [23, Ch. 2.1]. Consequently, we can rewrite Eq. (3) by multiplying with \mathbf{P} from the right:

$$\mathbf{X}_{:,1:j-1} \mathbf{P} = \mathbf{U}_{:,1:h} \mathbf{R}_{\text{ref}1:h,1:j-1} \mathbf{P} \implies \mathbf{U}_{:,1:h} = \mathbf{X}_{:,1:j-1} \mathbf{P}.$$

Thus, the first h columns in \mathbf{U} can be expressed using the first $j-1$ columns in \mathbf{X} as $\mathbf{u}_i = \sum_{k=1}^{j-1} \mathbf{x}_k P_{k,i}$, for $i \in \{1, \dots, h\}$. Inserting this into $\mathbf{x}_j = \sum_{i=1}^h \mathbf{u}_i r_{i,j}$ gives

$$\mathbf{x}_j = \sum_{i=1}^h \sum_{k=1}^{j-1} \mathbf{x}_k P_{k,i} r_{i,j} = \sum_{k=1}^{j-1} \mathbf{x}_k \sum_{i=1}^h P_{k,i} r_{i,j} = \sum_{k=1}^{j-1} c_{k,j} \mathbf{x}_k, \quad (4)$$

where $c_{k,j} = \sum_{i=1}^h P_{k,i} r_{i,j}$. Thus, for $h > 0$, we have shown that $\mathbf{r}_j \in \text{span}(\mathbf{r}_1, \dots, \mathbf{r}_{j-1}) \implies \mathbf{x}_j \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$.

The case $h = 0$ corresponds to a trivial case as it requires all columns of \mathbf{R}_{ref} with index smaller than j to consist of zeros. Consequently, all columns vectors of \mathbf{X} up to and including j must also consist of zeros, meaning that $\mathbf{x}_j \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$ trivially. \square

By this lemma, the linearly independent columns in \mathbf{X} correspond to the linearly independent columns in \mathbf{R}_{ref} , which can be easily identified by finding the pivot elements. However, the lemma requires a matrix decomposition of the form $\mathbf{X} = \mathbf{U} \mathbf{R}_{\text{ref}}$, where \mathbf{U} has a left-side inverse and \mathbf{R}_{ref} is in REF. The well-known QR decomposition $\mathbf{X} = \mathbf{Q} \mathbf{R}$ is a good starting point as the matrix $\mathbf{Q} \in \mathbb{R}^{n \times d}$ has orthonormal columns, so \mathbf{Q}^\top is its left-side inverse [12, Ch. 2.1]. Yet $\mathbf{R} \in \mathbb{R}^{d \times d}$ is only upper-triangular, not in REF.

Therefore, we put forward the QMR decomposition. The matrix \mathbf{R} is further decomposed as $\mathbf{R} = \mathbf{M} \mathbf{R}_{\text{ref}}$ using Gaussian elimination, such that $\mathbf{R}_{\text{ref}} \in \mathbb{R}^{d \times d}$ is in REF, and $\mathbf{M} \in \mathbb{R}^{d \times d}$ is an invertible matrix composed of inverted row reduction operations [9, Ch. 3.2]. That is, $\mathbf{M} = \mathbf{M}_k^{-1} \dots \mathbf{M}_2^{-1} \mathbf{M}_1^{-1}$, where each \mathbf{M}_k is one of three invertible operations: row swap, addition of two rows, and multiplication of a row with a nonzero value [9, Ch. 3.2.6]. Thus, we have $\mathbf{X} = \mathbf{Q} \mathbf{M} \mathbf{R}_{\text{ref}}$, with $\mathbf{U} = \mathbf{Q} \mathbf{M}$ having the left-side inverse $\mathbf{U}^\dagger = \mathbf{M}^{-1} \mathbf{Q}^\top$, fulfilling the criteria of Lemma 2.1. Computation of \mathbf{M} and \mathbf{R}_{ref} is implemented using outer product Gaussian elimination [9, Ch. 3.2.8], see lines 12 to 15 in Alg.1.

In practice, exact linear independence can be too strict, leading to more features being retained than desired. To address this, QMR-FS removes features with small reconstruction errors. Let \tilde{S} be the list of retained column indices, and let $\mathbf{X}_{:, \tilde{S}}$ denote the matrix of retained feature columns. We write $\mathbf{x}_j = \boldsymbol{\chi}_j + \boldsymbol{\delta}_j$, where $\boldsymbol{\chi}_j = \mathbf{X}_{:, \tilde{S}} \mathbf{c}$ denotes a reconstruction of \mathbf{x}_j using $\mathbf{X}_{:, \tilde{S}}$ with the coefficient vector \mathbf{c} , and $\boldsymbol{\delta}_j$ is a residual error vector. Then \mathbf{x}_j is removed if $\|\boldsymbol{\delta}_j\|_2 \leq \theta \|\mathbf{x}_j\|_2$, where $\theta \in [0, 1]$ is the relative tolerance. While \mathbf{c} and $\boldsymbol{\delta}_j$ can be computed via least-squares fitting, this approach is computationally expensive. Instead, we use a closed-form formula for $\boldsymbol{\delta}_j$, which provides an upper bound on the reconstruction error.

To derive the formula for $\boldsymbol{\delta}_j$, consider Alg. 1 at the start of iteration j , with $p = |\tilde{S}| + 1$. By the QMR decomposition, we have $\mathbf{x}_j = \mathbf{Q} \mathbf{M}_{:,1:p} \mathbf{R}_{1:p,j}$, which we can split into two terms:

$$\mathbf{x}_j = \mathbf{Q} \mathbf{M}_{:,1:p-1} \mathbf{R}_{1:p-1,j} + \mathbf{Q} \mathbf{M}_{:,p:j} \mathbf{R}_{p:j,j}. \quad (5)$$

Note that the first term is a linear combination of the first $p-1$ column vectors of the matrix $\mathbf{Q} \mathbf{M}$, while the second term linearly combines the next $j-p+1$ column vectors. We next show that these two terms correspond to $\boldsymbol{\chi}_j$ and $\boldsymbol{\delta}_j$ as defined above.

To see this, we show that the first term in Eq. (5) corresponds to a linear combination of the columns in $\mathbf{X}_{:, \tilde{S}}$. Let $\mathbf{R}' = \mathbf{R}_{1:p-1, \tilde{S}}$ denote the submatrix consisting of the first $p-1$ rows of \mathbf{R} and the columns with indices in \tilde{S} . By the algorithm's construction, \mathbf{R}' is in REF with full row rank, meaning \mathbf{R}' has the right-side inverse \mathbf{P} [23, Ch. 2.1]. Moreover, it holds that $\mathbf{X}_{:, \tilde{S}} = \mathbf{Q} \mathbf{M}_{:, 1:p-1} \mathbf{R}'$. To see why, first note that it is valid for $j = 0$. Next, observe that it is preserved through each iteration as \mathbf{M} , \mathbf{R} , p , and \tilde{S} are updated. Consider the two scenarios addressed by the if-statement on line 9. If the condition evaluates to False, Gaussian elimination proceeds normally, preserving the equation by construction. If the condition evaluates to True, the elements $\mathbf{R}_{p:j, j}$ are set to zero, ensuring that column j has no pivot elements in \mathbf{R}_{ref} . Importantly, since the elements in $\mathbf{R}_{p:j, j}$ are not used in subsequent iterations for updating \mathbf{M} , and j was not added to \tilde{S} , the equality $\mathbf{X}_{:, \tilde{S}} = \mathbf{Q} \mathbf{M}_{:, 1:p-1} \mathbf{R}'$ remains valid. Combining these facts, we see that $\mathbf{Q} \mathbf{M}_{:, 1:p-1} = \mathbf{X}_{:, \tilde{S}} \mathbf{P}$. Therefore, the first term in Eq. (5) equals χ_j , $\mathbf{Q} \mathbf{M}_{:, 1:p-1} \mathbf{R}_{1:p-1, j} = \mathbf{X}_{:, \tilde{S}} \mathbf{P} \mathbf{R}_{1:p-1, j} = \mathbf{X}_{:, \tilde{S}} \mathbf{c} = \chi_j$, meaning that the second term equals δ_j .

Three additional aspects of QMR-FS require attention. First, due to its greedy selection approach, QMR-FS is biased towards the initial feature ordering in \mathbf{X} , which results in features further to the left being more likely to be retained. Therefore, QMR-FS benefits from a good initial ordering of the features. If prior knowledge about feature importance is available, it should be used to establish the initial ordering. In other cases, we propose ordering features by their Shannon entropy [5, Ch. 2.1] (line 2). This heuristic is motivated by previously suggested entropy measures for UFS [7, 10, 28]. However, we acknowledge that improving the initial ordering is an important direction for future work on QMR-FS. Second, to account for the constant term in Eq. 1, we prepend $\mathbf{1}_n$ to \mathbf{X} (line 3 in Alg. 1). Finally, the time complexity of QMR-FS is $O(nd^2)$, dominated by line 8 and the QR decomposition [9, Ch. 5.2.2].

Algorithm 1: The QMR-FS algorithm with 1 based indexing. *Inputs:* Data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and tolerance threshold $\theta \in [0, 1]$.

```

1 def qmr_fs(X, θ):
2   X = sort_columns_by_entropy(X, 'descending')
3   X = [1n, X]      # Prepend column with constant value.
4   Q, R = qr_decomposition(X)
5   M = I_{d \times d}   # Initialize M as the identity matrix.
6   S̃ = list(), p = 1 # Initialize S̃ and p.
7   for j in range(1, R.shape[1]):
8     δ_j = Q M[:, p:j] R[p:j, j] # Formula provided by Eq. (5).
9     if ||δ_j||_2 ≤ θ ||x_j||_2:
10      R[p:j, j] = 0 # Ensure no pivot elements.
11    else:
12      # Gaussian elimination on column j from row p.
13      ρ = R[p+1:j, j] / R[p, j] # [(j-p-1) \times 1]
14      R[p+1:j, :] -= ρ \otimes R[p, :] # [(j-p-1) \times d]
15      M[:, p] += M[:, p+1:j] \cdot ρ # [d \times 1]
16      p += 1 # Increment pivot counter.
17      S̃.append(j) # Add feature j to selected set.
18      S̃.remove(1) # Remove constant value column.
19   return X[:, S̃]

```

3 Experiments

We conduct two sets of experiments to demonstrate the effectiveness of QMR-FS. First, we showcase its scalability by applying it to four large datasets: US Census (1990) [21], GitHub MUSAE [24], SNAP patents [17], and KDDCUP (1999) [27]. The results are shown in Tbl. 1. Notably, QMR-FS completes in just over 30 seconds on a CPU (16 vCPUs @ 2.2 GHz) and 20 seconds on a GPU (Nvidia L4) for KDDCUP, the largest dataset with over 1 billion elements, which is 10 times faster than recent scalable UFS methods [14, 15] with reported runtimes of 200s and 300s on a dataset with about half as many elements, $n = 630k$ and $d = 900$ [15].

In our second experiment, we demonstrate that QMR-FS performs on par with other UFS methods despite its greedy approach. Following the setup outlined in a recent review and benchmark paper [25], we evaluate the selected features using SVM classification [4] and K-means clustering [2]. For classification, we compute average accuracies using 5-fold cross-validation with five different seeds, and for clustering, we compute average normalized mutual information (NMI) [20, Ch. 16.3] using 25 seeds. We use five datasets from [25], all from the UCI ML repository [16], and add the popular Isolet dataset [3, 8, 19]. Dataset details are provided in Tbl. 2.

Following [25], we use SVD-entropy [28], LS [11], SPEC [34], USFSM [26], UDFS [31], and NDFS [18] as comparing UFS methods. To compensate for excluding methods in [25] without readily available implementations, we add the recent methods CNAFS [32] and FMIUFS [33], available in the Matlab UFS toolbox [1]. Each of these methods outputs a feature ranking and expects a specific proportion of features to be selected, which differs from QMR-FS, which automatically selects the number of features based on the tolerance threshold θ . To make results comparable, we extract eight different feature sets from each ranking produced by the baseline UFS methods, corresponding to 20%-90% of the features. For Isolet, we use a fixed number of features ranging from 20 to 100 instead of percentages, as is customary for this dataset [19, 29, 30]. For QMR-FS, we use multiple values of θ to obtain a spread over the number of preserved features.

Examining Fig. 1, we find that QMR-FS is among the top methods for classification; however, no single method excels across all settings. The average rank results in Tbl. 2 convey a similar story. QMR-FS achieves the highest classification ranking and an average clustering ranking within one standard deviation of the best method. Overall, the results support our claim that QMR-FS performs comparably to other UFS methods while offering superior scalability. Finally, the Isolet results in Fig. 1j are notable for QMR-FS's significant improvement in accuracy when increasing from 20 to 100 features, suggesting potential for enhancing the initial feature ordering.

Table 1: QMR-FS runtimes in seconds on four large datasets using $\theta = 0.1$, resulting in d_{fs} selected features.

| DATASET | NUM. INSTANCES AND DIMS. | | | RUNTIME (s) | |
|------------------|--------------------------|------|----------|-------------|-------------|
| | n | d | d_{fs} | CPU | GPU |
| US CENSUS (1990) | 2.46M | 68 | 66 | 4.16 ± 0.16 | 2.75 ± 0.11 |
| GITHub MUSAE | 37.7K | 4006 | 3799 | 28.5 ± 0.18 | 13.9 ± 0.45 |
| SNAP PATENTS | 2.92M | 269 | 259 | 26.4 ± 0.25 | 13.7 ± 0.06 |
| KDDCUP (1999) | 7.88M | 127 | 111 | 33.9 ± 0.12 | 20.2 ± 0.06 |

Table 2: Dataset details (left), and summary of benchmark results (right) using 40% and 60% kept features following [25] (#50 and #100 for Isolet). The average ranks and standard deviations for classification and clustering are computed over the 6 datasets. The time complexities are simplified under the assumption $n \geq d$, and * indicates methods which are iterated until convergence. Relative runtimes are displayed for Isolet, the largest dataset the baseline UFS methods scale to.

| DATASET | n | d | # CLASSES | | QMR-FS | SVD ENT. | LS | SPEC | USFSM | UDFS | NDFS | CNAFS | FMIUFS |
|---------------|------|-----|-----------|------------------------|---------------|---------------|---------------|---------------|----------------|---------------|---------------|---------------|----------------|
| AUTOMOBILE | 205 | 25 | 6 | CLSIF. AVG. RANK (40%) | 2.8 ± 1.8 | 5.5 ± 2.8 | 5.8 ± 2.6 | 8.3 ± 0.5 | 3.0 ± 2.2 | 5.5 ± 1.8 | 4.3 ± 1.9 | 5.0 ± 1.4 | 3.7 ± 2.6 |
| HEART-C | 303 | 13 | 5 | CLSIF. AVG. RANK (60%) | 2.5 ± 1.6 | 5.0 ± 1.8 | 5.7 ± 2.1 | 6.5 ± 2.6 | 5.8 ± 2.8 | 4.7 ± 2.8 | 3.2 ± 1.7 | 5.3 ± 3.1 | 5.2 ± 3.1 |
| HEART-STATLOG | 270 | 13 | 2 | CLSTR. AVG. RANK (40%) | 3.0 ± 2.2 | 6.7 ± 2.1 | 6.5 ± 2.6 | 7.0 ± 1.7 | 3.7 ± 2.1 | 6.0 ± 2.8 | 3.3 ± 2.3 | 5.5 ± 1.6 | 2.8 ± 1.3 |
| SONAR | 208 | 60 | 2 | CLSTR. AVG. RANK (60%) | 5.5 ± 2.9 | 4.3 ± 2.0 | 4.8 ± 2.4 | 7.2 ± 1.3 | 4.8 ± 2.5 | 5.5 ± 2.6 | 3.0 ± 2.7 | 4.2 ± 3.4 | 4.8 ± 2.3 |
| WINE | 178 | 13 | 3 | TIME COMPLEXITY | $O(nd^2)$ | $O(nd^3)$ | $O(n^2d)$ | $O(n^2d)$ | $O(n^3d)$ | $O(n^2d)^*$ | $O(n^2d)^*$ | $O(n^2d)^*$ | $O(n^2d)$ |
| ISOLET | 1560 | 617 | 26 | RUNTIME (ISOLET) | 137MS | $\times 310$ | $\times 1.2$ | $\times 16$ | $\times 22764$ | $\times 11$ | $\times 44$ | $\times 2950$ | $\times 13635$ |

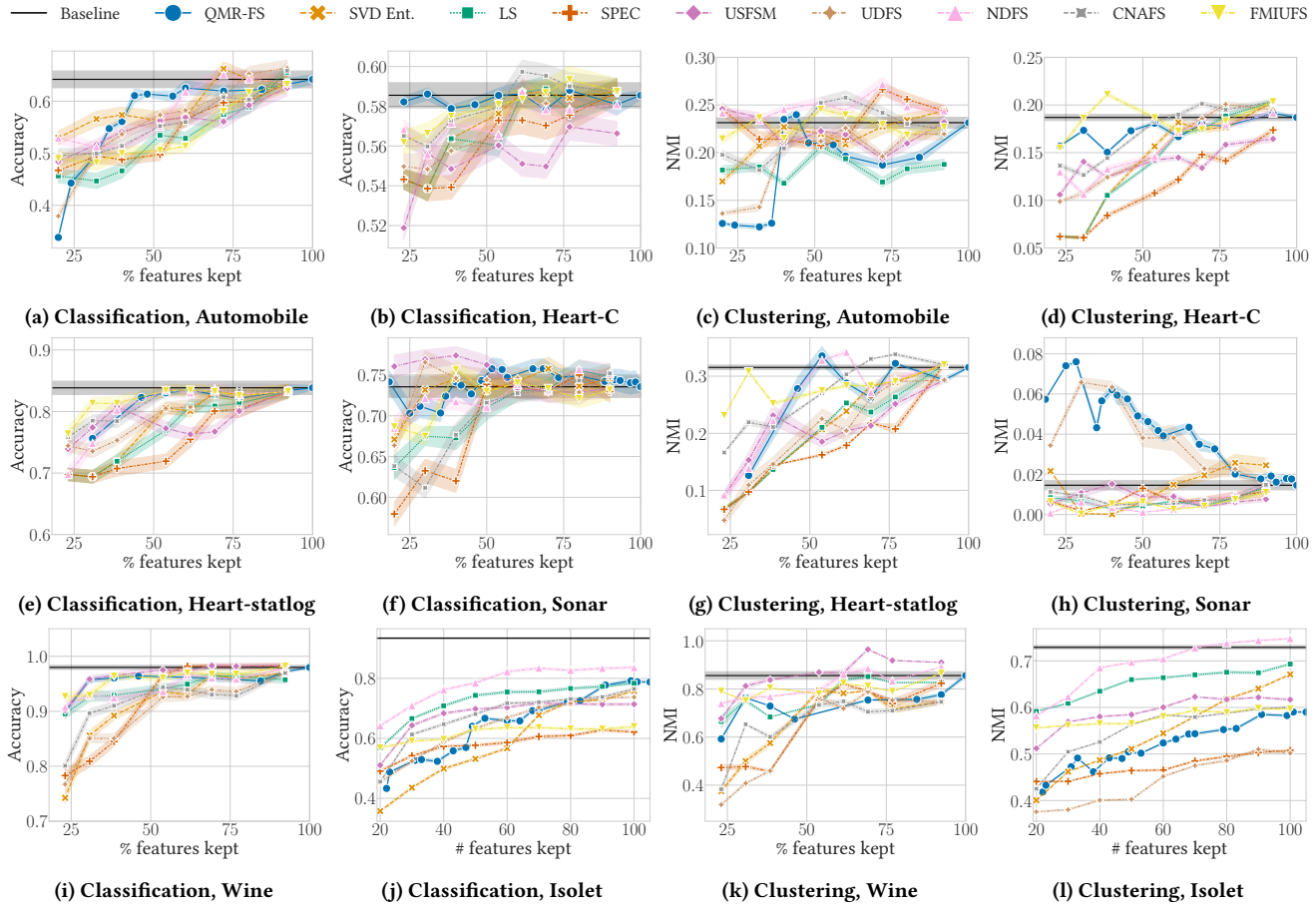


Figure 1: Feature selection benchmark results. Each marker shows average SVM classification accuracy (left) or K-means NMI (right). The shaded areas show the standard error. The solid black lines show the score when 100% of features are used.

4 Conclusion and Future Work

In this paper, we have presented QMR-FS, a scalable unsupervised feature selection (UFS) method, demonstrating its ability to achieve classification and clustering accuracies comparable to existing UFS methods while performing efficiently on large datasets. Despite these strengths, there are several avenues for further improvement. Enhancing the initial feature ordering, to which QMR-FS is biased, holds significant potential. This could be achieved by refining the feature ordering heuristic or by running QMR-FS with multiple

random orderings and selecting the best order based on total reconstruction error. Given the speed of QMR-FS, this approach would not impose substantial computational overhead.

Another research direction involves reducing the bias associated with feature ordering. For example, QMR-FS could be run iteratively, altering the feature order between iterations to reveal feature reconstruction possibilities not considered with a single ordering. Additionally, the differentiable nature of the QMR decomposition opens up exciting possibilities for applications in the deep learning domain, which we are eager to explore.

References

- [1] Farhad Abedinzadeh Torghabeh, Yeganeh Modaresnia, and Seyyed Abed Hosseini. 2023. Auto-UFSTool: An Automatic Unsupervised Feature Selection Toolbox for MATLAB. *Journal of AI and Data Mining* 11, 4 (2023), 517–524. <https://doi.org/10.22044/jadm.2023.12820.2434>
- [2] David Arthur and Sergei Vassilvitskii. 2007. K-Means++: The Advantages of Careful Seeding. In *SODA '07*. ACM, 1027–1035. <https://dl.acm.org/doi/10.5555/1283383.1283494>
- [3] Deng Cai, Xiaofei He, and Jiawei Han. 2011. Speed up kernel discriminant analysis. *The VLDB Journal* 20 (2011), 21–33. <https://doi.org/10.1007/s00778-010-0189-3>
- [4] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20 (1995), 273–297. <https://doi.org/10.1007/BF00994018>
- [5] Thomas M. Cover and Joy A. Thomas. 2005. *Entropy, Relative Entropy, and Mutual Information*. John Wiley & Sons, Ltd. 13–55 pages. <https://doi.org/10.1002/047174882X.ch2>
- [6] S.K. Das. 1971. Feature Selection with a Linear Dependence Measure. *IEEE Trans. Comput.* C-20, 9 (1971), 1106–1109. <https://doi.org/10.1109/T-C.1971.223412>
- [7] M. Dash, H. Liu, and J. Yao. 1997. Dimensionality reduction of unsupervised data. In *Proceedings Ninth IEEE International Conference on Tools with Artificial Intelligence*. IEEE, 532–539. <https://doi.org/10.1109/TAI.1997.632300>
- [8] Mark Fandy and Ronald Cole. 1990. Spoken Letter Recognition. In *Advances in Neural Information Processing Systems*, Vol. 3. Morgan-Kaufmann, 220–226. <https://proceedings.neurips.cc/paper/1990>
- [9] Gene H Golub and Charles F Van Loan. 2013. *Matrix computations* (fourth ed.). JHU press. <https://doi.org/10.5555/944919.944968>
- [10] Isabelle Guyon and André Elisseeff. 2003. An introduction to variable and feature selection. *Journal of machine learning research* 3, Mar (2003), 1157–1182. <https://dl.acm.org/doi/10.5555/944919.944968>
- [11] Xiaofei He, Deng Cai, and Partha Niyogi. 2005. Laplacian score for feature selection. In *NIPS'05*. MIT Press, 507–514. <https://dl.acm.org/doi/abs/10.5555/2976248.2976312>
- [12] Roger A. Horn and Charles R. Johnson. 2012. *Matrix Analysis* (2 ed.). Cambridge University Press. <https://doi.org/10.1017/CBO9781139020411>
- [13] Michael E. Houle, Hans-Peter Kriegel, Peer Kröger, Erich Schubert, and Arthur Zimek. 2010. Can Shared-Neighbor Distances Defeat the Curse of Dimensionality?. In *Scientific and Statistical Database Management*. Springer Berlin Heidelberg, 482–500. https://doi.org/10.1007/978-3-642-13818-8_34
- [14] Haojie Hu, Rong Wang, Feiping Nie, Xiaojun Yang, and Weizhong Yu. 2018. Fast unsupervised feature selection with anchor graph and $\ell_{2,1}$ -norm regularization. *Multimedia Tools Appl.* 77, 17 (sep 2018), 22099–22113. <https://doi.org/10.1007/s11042-017-5582-0>
- [15] Haojie Hu, Rong Wang, Xiaojun Yang, and Feiping Nie. 2019. Scalable and Flexible Unsupervised Feature Selection. *Neural Computation* 31, 3 (2019), 517–537. https://doi.org/10.1162/neco_a_01163
- [16] Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. 2024. The UCI Machine Learning Repository. <https://archive.ics.uci.edu>
- [17] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. 2005. SNAP Patents. Stanford Large Network Dataset Collection. <https://snap.stanford.edu/data/cit-Patents.html>
- [18] Zechao Li, Yi Yang, Jing Liu, Xiaofang Zhou, and Hanqing Lu. 2012. Unsupervised feature selection using nonnegative spectral analysis. In *AAAI'12*. AAAI Press, 1026–1032. <https://doi.org/10.1609/aaai.v26i1.8289>
- [19] Yanfang Liu, Dongyi Ye, Wenbin Li, Huihui Wang, and Yang Gao. 2020. Robust neighborhood embedding for unsupervised feature selection. *Knowledge-Based Systems* 193 (2020), 105462. <https://doi.org/10.1016/j.knsys.2019.105462>
- [20] Christopher D Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to information retrieval*. Cambridge University Press. <https://nlp.stanford.edu/IR-book/information-retrieval-book.html>
- [21] Chris Meek, Bo Thiesson, and David Heckerman. 2001. US Census Data (1990). UCI Machine Learning Repository. <https://doi.org/10.24432/C5VP42>
- [22] Pabitra Mitra, Chivukula Anjaneya Murthy, and Sankar K. Pal. 2002. Unsupervised feature selection using feature similarity. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24, 3 (2002), 301–312. <https://doi.org/10.1109/34.990133>
- [23] C. Radhakrishna Rao. 1971. *Generalized inverse of matrices and its applications*. Wiley-Interscience. This book is not available online. Readers who cannot access a physical copy can also look in this online document for the same statement: MIT course document.
- [24] UCI Machine Learning Repository. 2019. GitHub MUSAE. UCI Machine Learning Repository. <https://doi.org/10.24432/C5Z02B>
- [25] Saul Solorio-Fernández, Ariel Carrasco-Ochoa, and José Fco Martínez-Trinidad. 2020. A review of unsupervised feature selection methods. *Artificial Intelligence Review* 53, 2 (2020), 907–948. <https://doi.org/10.1007/s10462-019-09682-y>
- [26] Saul Solorio-Fernández, José Fco. Martínez-Trinidad, and J. Ariel Carrasco-Ochoa. 2017. A new Unsupervised Spectral Feature Selection Method for mixed data: A filter approach. *Pattern Recognition* 72 (2017), 314–326. <https://doi.org/10.1016/j.patcog.2017.07.020>
- [27] Salvatore Stolfo, Wei Fan, Wenke Lee, Andreas Prodromidis, and Philip Chan. 1999. KDD Cup 1999 Data. UCI Machine Learning Repository. <https://doi.org/10.24432/C51C7N>
- [28] Roy Varshavsky, Assaf Gottlieb, Michal Linial, and David Horn. 2006. Novel Unsupervised Feature Filtering of Biological Data. *Bioinformatics* 22, 14 (07 2006), e507–e513. <https://doi.org/10.1093/bioinformatics/btl214>
- [29] Shiping Wang, Witold Pedrycz, Qingxin Zhu, and William Zhu. 2015. Subspace learning for unsupervised feature selection via matrix factorization. *Pattern Recognition* 48, 1 (2015), 10–19. <https://doi.org/10.1016/j.patcog.2014.08.004>
- [30] Shiping Wang, Witold Pedrycz, Qingxin Zhu, and William Zhu. 2015. Unsupervised feature selection via maximum projection and minimum redundancy. *Knowledge-Based Systems* 75 (2015), 19–29. <https://doi.org/10.1016/j.knsys.2014.11.008>
- [31] Yi Yang, Heng Tao Shen, Zhigang Ma, Zi Huang, and Xiaofang Zhou. 2011. $\ell_{2,1}$ -norm regularized discriminative feature selection for unsupervised learning. In *IJCAI'11*. AAAI Press, 1589–1594. <https://dl.acm.org/doi/10.5555/2283516.2283660>
- [32] Aihong Yuan, Mengbo You, Dongjian He, and Xuelong Li. 2022. Convex Non-Negative Matrix Factorization With Adaptive Graph for Unsupervised Feature Selection. *IEEE Transactions on Cybernetics* 52, 6 (2022), 5522–5534. <https://doi.org/10.1109/TCYB.2020.3034462>
- [33] Zhong Yuan, Hongmei Chen, Pengfei Zhang, Jihong Wan, and Tianrui Li. 2022. A Novel Unsupervised Approach to Heterogeneous Feature Selection Based on Fuzzy Mutual Information. *IEEE Transactions on Fuzzy Systems* 30, 9 (2022), 3395–3409. <https://doi.org/10.1109/TFUZZ.2021.3114734>
- [34] Zheng Zhao and Huan Liu. 2007. Spectral feature selection for supervised and unsupervised learning. In *ICML'07*. ACM, 1151–1157. <https://doi.org/10.1145/1273496.1273641>