Scalable Unsupervised Feature Selection with Reconstruction Error Guarantees via QMR Decomposition



KTH Royal Institute of Technology Division of Robotics, Perception, and Learning

Abstract

Unsupervised feature selection (UFS) methods are valuable for eliminating redundant features without using class labels, but they often struggle with large datasets. To overcome this, we present QMR-FS, a scalable, greedy forward filtering approach that selects linearly independent features while ensuring any excluded features can be reconstructed within a specified tolerance. QMR-FS is based on QMR matrix decomposition, an extension of QR decomposition, resulting in linear complexity relative to the number of samples. This facilitates parallelized computation on both CPU and GPU, and our model implementation runs in a matter of seconds on datasets with up to 1 billion elements, while offering classification and clustering performance comparable to other UFS methods.

Problem: The UFS task is to select a subset of the columns of a data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, while preserving or improving downstream task accuracy. We desire a method which scales to large *n*.

QMR-FS is a greedy approach



. 1.4

The lemma requires a matrix decomposition $X=UR_{ref}$ where U has a left-side inverse, and R_{ref} is in row echelon form. The QR decomposition is a good start, but insufficient since **R** is only upper triangular. Here is an example showcasing this limitation.

SEB Group

SEB

	[1	1	2	1	1]		-0.45	0.89	0.	0.	0.]		-2.24	-0.45	-0.89	-1.34	-1.79
	1	0	0	1	2		-0.45	-0.22	0.87	0.	0.		0.	0.89	1.79	0.45	0.22
$\mathbf{X} =$	1	0	0	1	1	$\mathbf{Q} =$	-0.45	-0.22	-0.29	0.82	0.	$\mathbf{R} =$	0.	0.	0.	0.58	1.44
	1	0	0	0	0		-0.45	-0.22	-0.29	-0.41	-0.71		0.	0.	0.	0.82	0.82
	1	0	0	0	0		-0.45	-0.22	-0.29	-0.41	0.71		0.	0.	0.	0.	0.

based on linear independence. $\mathbf{x}_j = \sum_{k=1}^{j} a_k \mathbf{x}_k + b \mathbf{1}_n$

Iterating over the columns in **X** from left to right, we drop columns which can be expressed as a linear combination of previously seen columns. Specifically, column *j* is dropped if there exists coefficients a_k and *b* that fulfills the formula above.

Although this could be done by performing d least-squares fits, it would be inefficient, with a time complexity of $O(nd^3)$.

Instead, QMR-FS archives time complexity *O(nd²)* by utilizing the following lemma, which we prove in our paper.

LEMMA 2.1. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{R}_{ref} \in \mathbb{R}^{d \times d}$ satisfy $\mathbf{X} = \mathbf{U}\mathbf{R}_{ref}$ and $\mathbf{R}_{ref} = \mathbf{U}^{\dagger} \mathbf{X}$, 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 jth column in \mathbf{X} , and $\mathbf{r}_{j} \in \mathbb{R}^{d}$ is the jth column in \mathbf{R}_{ref} . As **R** is not in row echelon form, it is not clear that columns 4 and 5 are both linearly independent. Therefore, **R** must be further decomposed using Gaussian elimination. This results in the QMR decomposition: **X=QMR**_{ref} :

	[1.	0.	0.	0.	0.]		-2.2	-0.45	-0.89	-1.34	-1.79
	0.	1.	0.	0.	0.		0.	0.89	1.79	0.45	0.22
$\mathbf{M} =$	0.	0.	1.	0.	0.	$\mathbf{R}_{ ext{ref}} =$	0.	0.	0.	0.58	1.44
	0.	0.	1.41	1.	0.		0.	0.	0.	0.	-1.23
	0.	0.	0.	0.	1.		0.	0.	0.	0.	0.

Gaussian elimination is invertible, meaning that **QM** has the left inverse matrix $\mathbf{M}^{-1}\mathbf{Q}^{\mathsf{T}}$. Thus, by the lemma, the linearly independent columns in **X** are the same as those in \mathbf{R}_{ref} , which are identified by the pivot elements. In the example, column 3 can be removed.

Thresholding: 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. During the Gaussian elimination, we can efficiently compute an upper bound on the reconstruction error without expensive least-square fitting. See our paper for more details.

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

2	NUM. IN	ISTANCE	ES AND DIMS.	RUNTIME (S)			
Dataset	n	d	$d_{ m fs}$	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		

Table 2: Summary of benchmark results using 40% and 60% kept features. The average ranks and standard deviations for classification and clustering are computed over the 6 datasets. The time complexities are simplified under the assumption $n \ge d$, and * indicates methods which are iterated until convergence. Relative runtimes are displayed for Isolet, the largest dataset the baseline UFS methods scale to.

	QMR-FS	SVD Ent.	LS	SPEC	USFSM	UDFS	NDFS	CNAFS	FMIUFS
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
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
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
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





