

# A Theoretical Analysis for CUR Decomposition based Active Learning and Feature Selection

**Zhong Chen**

ZHONG.CHEN@SIU.EDU

*School of Computing, Southern Illinois University, Carbondale, IL 62901, United States*

**Chen Zhao**

CHEN\_ZHAO@BAYLOR.EDU

*Department of Computer Science, Baylor University, Waco, TX 76798, United States*

**Yi He**

YIHE@WM.EDU

*Department of Data Science, William & Mary, Williamsburg, VA 23185, United States*

## Abstract

This study comprehensively investigated CUR decomposition-based active learning and feature selection from an optimization perspective, especially provided the complexity analysis.

## 1. Introduction

The primary goal of active learning [4, 12, 15, 21, 24] and feature selection [6, 9, 11, 13, 25, 26, 28, 30] is to address a fundamental challenge in machine learning and online learning: the high cost of labeling data and the curse of dimensionality. Modern datasets often contain a vast number of both instances (rows) and features (columns), making full annotation prohibitively expensive and models computationally heavy and prone to overfitting. The paper introduces a unified solution that tackles both problems simultaneously by leveraging the principles of matrix decomposition, specifically the CUR matrix decomposition [3, 5, 12, 14, 23], to select the most informative data points to label (active learning) and the most relevant features (feature selection).

CUR decomposition is chosen over other techniques such as PCA [22] or SVD [1] because it provides a low-rank approximation of the original data matrix using actual rows and columns from the dataset itself. This crucial property ensures that the selected features and instances are interpretable and maintain their original meaning, as opposed to the transformed, abstract components produced by PCA. The method works by minimizing the residual of CUR with sparse row and column constraints, which identify the rows (feature) and columns (instances) that exert the most influence on the structure of the data matrix. These high-leverage points are deemed the most representative and informative.

By applying CUR decomposition, the framework efficiently selects a small subset of instances for an expert to label, drastically reducing annotation costs in active learning. Concurrently, it selects a subset of features that best capture the data's variance, improving model efficiency and generalization. This creates a powerful synergy: the model is trained on maximally informative data points described by the most relevant features, leading to a more robust, interpretable, and computationally efficient machine learning pipeline compared to handling each problem separately.

## 2. Active Learning and Feature Selection via CUR Decomposition

From an algorithmic perspective, the matrices  $\mathbf{C}$ ,  $\mathbf{U}$ , and  $\mathbf{R}$  can be obtained by minimizing the approximation error  $\|\mathbf{W} - \mathbf{CUR}\|_F^2$ . Here we make a key observation that the above definition is closely related to the problem of simultaneous sample and feature selection. More specifically, the matrix  $\mathbf{UR}$  can be regarded as a reconstruction coefficient matrix, and  $\mathbf{C}$  denotes the selected  $m$  samples, thus minimizing means that the total reconstruction error  $\|\mathbf{W} - \mathbf{CUR}\|_F^2$  is minimized, which can make the data points in  $\mathbf{C}$  be the most representative. The reconstruction coefficients  $\mathbf{UR}$  are related to an  $r$ -dimensional feature subset of the dataset. The reconstruction coefficients of each reconstructed data point  $\mathbf{w}_i$  ( $i = 1, 2, \dots, n$ ) are formed by a linear combination of its  $r$  features. In the meantime, the matrix  $\mathbf{CU}$  can also be regarded as a reconstruction coefficient matrix, and  $\mathbf{R}$  is the new low-dimensional representation of  $\mathbf{W}$ , so minimizing  $\|\mathbf{W} - \mathbf{CUR}\|_F^2$  also indicates that the selected  $r$  features can represent the whole dataset most precisely. The construction of the coefficient matrix  $\mathbf{CU}$  depends on a sample subset of  $\mathbf{W}$ . Clearly, active learning and feature selection can be conducted simultaneously in such a joint framework via CUR factorization.

Let  $\mathbf{p} = [p_1, p_2, \dots, p_n]^T \in \{0, 1\}^n$  and  $\mathbf{q} = [q_1, q_2, \dots, q_d]^T \in \{0, 1\}^d$  denote two indicator variables to represent whether a sample and a feature is selected or not, respectively. Specifically,  $p_i = 1$  (or 0) ( $i = 1, 2, \dots, n$ ) indicates that the  $i$ -th sample is selected (or not), and  $q_j = 1$  (or 0) ( $j = 1, 2, \dots, d$ ) means that the  $j$ -th feature is selected (or not). Then, minimizing  $\|\mathbf{W} - \mathbf{CUR}\|_F^2$  can be rewritten as

$$\begin{aligned} \min_{\mathbf{p} \in \{0,1\}^n, \mathbf{V} \in \mathbb{R}^{n \times d}, \mathbf{q} \in \{0,1\}^d} & \|\mathbf{W} - \mathbf{W} \text{diag}(\mathbf{p}) \mathbf{V} \text{diag}(\mathbf{q}) \mathbf{W}\|_F^2 \\ \text{s.t. } & \mathbf{1}_n^T \mathbf{p} = m, \mathbf{p} \in \{0, 1\}^n \\ & \mathbf{1}_d^T \mathbf{q} = r, \mathbf{q} \in \{0, 1\}^d \end{aligned} \quad (1)$$

where  $\mathbf{V} \in \mathbb{R}^{n \times d}$ ,  $\text{diag}(\mathbf{p}) = \text{diag}\{p_1, p_2, \dots, p_n\} \in \mathbb{R}^{n \times n}$  is a diagonal matrix with  $\mathbf{p}$  on its diagonal,  $\text{diag}(\mathbf{q}) = \text{diag}\{q_1, q_2, \dots, q_d\} \in \mathbb{R}^{d \times d}$  is a diagonal matrix with  $\mathbf{q}$  on its diagonal,  $\mathbf{1}_n = [1, 1, \dots, 1]^T \in \mathbb{R}^n$  is an  $n$ -dimensional vector with all components being 1, and  $\mathbf{1}_d = [1, 1, \dots, 1]^T \in \mathbb{R}^d$  is a  $d$ -dimensional vector with all components being 1. The term  $\mathbf{W} \text{diag}(\mathbf{p})$  in Eq. (1) aims to make  $m$  columns of  $\mathbf{W}$  unchanged, and resets the rest  $(n - m)$  columns to zero vectors. While the term  $\text{diag}(\mathbf{q}) \mathbf{W}$  in Eq. (1) tends to keep  $r$  rows of  $\mathbf{W}$  unchanged, and resets the rest  $(d - r)$  row to zero vectors.

## 3. A Convex Formulation

The objective function (1) is hard to be solved directly, since it is an NP-hard problem. After a careful observation to (1), we find that we can utilize the  $\ell_{2,0}$  mixed norm of a matrix to reduce the number of the parameters. Defining  $\mathbf{X} = \text{diag}(\mathbf{p}) \mathbf{V} \text{diag}(\mathbf{q}) \in \mathbb{R}^{n \times d}$ , we can rewritten Eq. (1) as

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times d}} & \|\mathbf{W} - \mathbf{W} \mathbf{X} \mathbf{W}\|_F^2 \\ \text{s.t. } & \|\mathbf{X}\|_{2,0} = m \\ & \|\mathbf{X}^T\|_{2,0} = r \end{aligned} \quad (2)$$

where the quasi-norm  $\ell_{2,0}$  norm of a matrix  $\mathbf{X}$  is defined as the number of the non-zero rows of  $\mathbf{X}$ , denoted by  $\|\mathbf{X}\|_{2,0}$ . Based on Eq. (2), we propose to optimize the following objective function:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times d}} \|\mathbf{W} - \mathbf{W}\mathbf{X}\mathbf{W}\|_F^2 + \alpha \|\mathbf{X}\|_{2,0} + \beta \|\mathbf{X}^T\|_{2,0} \quad (3)$$

where  $\alpha > 0$  and  $\beta > 0$  are two regularization parameters.

However, Eq. (3) is still an NP-hard problem due to the  $\ell_{2,0}$  mixed norm of a matrix. Fortunately, there exists theoretical progress that  $\|\mathbf{X}\|_{2,1}$  is the minimum convex hull of  $\|\mathbf{X}\|_{2,0}$ . The result of minimizing  $\|\mathbf{X}\|_{2,1}$  is the same as that of minimizing  $\|\mathbf{X}\|_{2,0}$ , as long as  $\mathbf{X}$  is row-sparse enough. Therefore, Eq. (3) can be relaxed to the following convex optimization problem:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times d}} \|\mathbf{W} - \mathbf{W}\mathbf{X}\mathbf{W}\|_F^2 + \alpha \|\mathbf{X}\|_{2,1} + \beta \|\mathbf{X}^T\|_{2,1} \quad (4)$$

where  $\|\mathbf{X}\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^d x_{i,j}^2} = \sum_{i=1}^n \|\mathbf{x}_i\|_2$ ,  $x_{i,j}$  is the  $i$ -th row  $j$ -th column element of  $\mathbf{X}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$  is the  $i$ -th row of  $\mathbf{X}$ .

#### 4. Optimization Algorithm

To solve Eq. (4), we first introduce two variables  $\mathbf{Y}$  and  $\mathbf{Z}$ , to convert to the following equivalent objective function:

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{Y} \in \mathbb{R}^{n \times d}, \mathbf{Z} \in \mathbb{R}^{d \times n}} & \|\mathbf{W} - \mathbf{W}\mathbf{X}\mathbf{W}\|_F^2 + \alpha \|\mathbf{Y}\|_{2,1} + \beta \|\mathbf{Z}\|_{2,1} \\ \text{s.t. } & \mathbf{Y} = \mathbf{X} \\ & \mathbf{Z} = \mathbf{X}^T \end{aligned} \quad (5)$$

The augmented Lagrange function of Eq. (5) is:

$$\begin{aligned} \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{\Lambda}_1, \mathbf{\Lambda}_2) &= \|\mathbf{W} - \mathbf{W}\mathbf{X}\mathbf{W}\|_F^2 + \alpha \|\mathbf{Y}\|_{2,1} + \beta \|\mathbf{Z}\|_{2,1} \\ &+ \langle \mathbf{\Lambda}_1, \mathbf{X} - \mathbf{Y} \rangle + \langle \mathbf{\Lambda}_2, \mathbf{X}^T - \mathbf{Z} \rangle \\ &+ \frac{\rho_1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 + \frac{\rho_2}{2} \|\mathbf{X}^T - \mathbf{Z}\|_F^2 \end{aligned} \quad (6)$$

where  $\mathbf{\Lambda}_1 \in \mathbb{R}^{n \times d}$  and  $\mathbf{\Lambda}_2 \in \mathbb{R}^{d \times n}$  are the Lagrange multipliers,  $\langle \mathbf{\Lambda}_1, \mathbf{X} - \mathbf{Y} \rangle = \text{Trace}(\mathbf{\Lambda}_1^T (\mathbf{X} - \mathbf{Y}))$ ,  $\text{Trace}(\cdot)$  is the trace norm of a square matrix, which is the sum of diagonal entries of a square matrix,  $\rho_1$  and  $\rho_2$  are the constraint violation penalty parameters. From the augmented Lagrangian function, we find that the subproblems about  $\mathbf{Y}$  and  $\mathbf{Z}$  are fully separable, as a result we can introduce the classical two-block ADMM [16], while considering  $\mathbf{X}$  and  $(\mathbf{Y}, \mathbf{Z})$  as two-block variables. In an ADMM-type algorithm [17], the basic Gauss-Seidel structure [2] in  $(t+1)$ -th iteration is

$$\begin{aligned} \mathbf{X}^{(t+1)} &= \arg\min_{\mathbf{X} \in \mathbb{R}^{n \times d}} \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}, \mathbf{Y}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)}) \\ \mathbf{Y}^{(t+1)} &= \arg\min_{\mathbf{Y} \in \mathbb{R}^{n \times d}} \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}^{(t+1)}, \mathbf{Y}, \mathbf{Z}^{(t)}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)}) \\ \mathbf{Z}^{(t+1)} &= \arg\min_{\mathbf{Z} \in \mathbb{R}^{d \times n}} \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}^{(t+1)}, \mathbf{Y}^{(t+1)}, \mathbf{Z}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)}) \\ \mathbf{\Lambda}_1^{(t+1)} &= \mathbf{\Lambda}_1^{(t)} + \rho_1 (\mathbf{X}^{(t+1)} - \mathbf{Y}^{(t+1)}) \\ \mathbf{\Lambda}_2^{(t+1)} &= \mathbf{\Lambda}_2^{(t)} + \rho_2 ((\mathbf{X}^{(t+1)})^T - \mathbf{Z}^{(t+1)}) \end{aligned} \quad (7)$$

Next, we will introduce how to solve these subproblems in detail.

(1) Compute the subproblem about  $\mathbf{X}^{(t+1)}$ :

When the other variables are fixed with the former iteration result  $(\mathbf{Y}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)})$ , the subproblem about  $\mathbf{X}^{(t+1)}$  is as

$$\begin{aligned} \mathbf{X}^{(t+1)} &= \operatorname{argmin}_{\mathbf{X} \in \mathbb{R}^{n \times d}} \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}, \mathbf{Y}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)}) \\ &= \operatorname{argmin}_{\mathbf{X} \in \mathbb{R}^{n \times d}} \{ \|\mathbf{W} - \mathbf{W}\mathbf{X}\mathbf{W}\|_F^2 \\ &\quad + \frac{\rho_1}{2} \|\mathbf{X} - \mathbf{Y}^{(t)} + \frac{1}{\rho_1} \mathbf{\Lambda}_1^{(t)}\|_F^2 \\ &\quad + \frac{\rho_2}{2} \|\mathbf{X}^T - \mathbf{Z}^{(t)} + \frac{1}{\rho_2} \mathbf{\Lambda}_2^{(t)}\|_F^2 \} \end{aligned} \quad (8)$$

The necessary optimality condition further follows as

$$\frac{\partial \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}, \mathbf{Y}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)})}{\partial \mathbf{X}} = 0 \quad (9)$$

This implies

$$\begin{aligned} &2\mathbf{W}^T \mathbf{W} \mathbf{X} \mathbf{W} \mathbf{W}^T + (\rho_1 + \rho_2) \mathbf{X} \\ &= 2\mathbf{W}^T \mathbf{W} \mathbf{W}^T + \rho_1 (\mathbf{Y}^{(t)} - \frac{1}{\rho_1} \mathbf{\Lambda}_1^{(t)}) + \rho_2 (\mathbf{Z}^{(t)} - \frac{1}{\rho_2} \mathbf{\Lambda}_2^{(t)})^T \end{aligned} \quad (10)$$

For writing conveniently, let  $\mathbf{M} = \mathbf{W}^T \mathbf{W} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{N} = \mathbf{W} \mathbf{W}^T \in \mathbb{R}^{d \times d}$ , and  $\mathbf{H} = 2\mathbf{W}^T \mathbf{W} \mathbf{W}^T + \rho_1 (\mathbf{Y}^{(t)} - \frac{1}{\rho_1} \mathbf{\Lambda}_1^{(t)}) + \rho_2 (\mathbf{Z}^{(t)} - \frac{1}{\rho_2} \mathbf{\Lambda}_2^{(t)})^T \in \mathbb{R}^{n \times d}$ , then, the equation above becomes

$$2\mathbf{M} \mathbf{X} \mathbf{N} + (\rho_1 + \rho_2) \mathbf{X} = \mathbf{H} \quad (11)$$

Since  $\mathbf{M} = \mathbf{W}^T \mathbf{W} \in \mathbb{R}^{n \times n}$  and  $\mathbf{N} = \mathbf{W} \mathbf{W}^T \in \mathbb{R}^{d \times d}$  are both positive semi-definite and symmetric, we can perform eigenvalue decomposition with all non-negative eigenvalues, obtaining

$$\begin{aligned} \mathbf{M} &= \mathbf{W}^T \mathbf{W} = \mathbf{P} \mathbf{\Phi} \mathbf{P}^T \in \mathbb{R}^{n \times n} \\ \mathbf{N} &= \mathbf{W} \mathbf{W}^T = \mathbf{Q} \mathbf{\Psi} \mathbf{Q}^T \in \mathbb{R}^{d \times d} \end{aligned} \quad (12)$$

where  $\mathbf{P} \in \mathbb{R}^{n \times n}$  and  $\mathbf{Q} \in \mathbb{R}^{d \times d}$  are two orthogonal matrices with all eigenvectors,  $\mathbf{P}^T \mathbf{P} = \mathbf{I}_n \in \mathbb{R}^{n \times n}$ ,  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_d \in \mathbb{R}^{d \times d}$ , and  $\mathbf{\Phi} \in \mathbb{R}^{n \times n}$  and  $\mathbf{\Psi} \in \mathbb{R}^{d \times d}$  are two diagonal matrices with all non-negative eigenvalues.

Plugging Eq. (12) into Eq. (11), we obtain

$$2\mathbf{P} \mathbf{\Phi} \mathbf{P}^T \mathbf{X} \mathbf{Q} \mathbf{\Psi} \mathbf{Q}^T + (\rho_1 + \rho_2) \mathbf{X} = \mathbf{H} \quad (13)$$

That is

$$2\mathbf{P}^T \mathbf{P} \mathbf{\Phi} \mathbf{P}^T \mathbf{X} \mathbf{Q} \mathbf{\Psi} \mathbf{Q}^T \mathbf{Q} + (\rho_1 + \rho_2) \mathbf{P}^T \mathbf{X} \mathbf{Q} = \mathbf{P}^T \mathbf{H} \mathbf{Q} \quad (14)$$

That is

$$2\mathbf{\Phi} \mathbf{P}^T \mathbf{X} \mathbf{Q} \mathbf{\Psi} + (\rho_1 + \rho_2) \mathbf{P}^T \mathbf{X} \mathbf{Q} = \mathbf{P}^T \mathbf{H} \mathbf{Q} \quad (15)$$

Let  $P^T X Q = D \in \mathbb{R}^{n \times d}$ , then Eq. (15) becomes

$$2\Phi D \Psi + (\rho_1 + \rho_2)D = P^T H Q \quad (16)$$

Thus, the  $i$ -th row and  $j$ -th column element  $D_{i,j}$  of the matrix  $D$  can be obtained by the following closed-form solution

$$D_{i,j} = \frac{(P^T H Q)_{i,j}}{2(\Phi)_{i,i}(\Psi)_{j,j} + (\rho_1 + \rho_2)} \quad (17)$$

where  $i = 1, 2, \dots, n$ ,  $j = 1, 2, \dots, d$ ,  $(P^T H Q)_{i,j}$  is the  $i$ -th row and  $j$ -th column element of the matrix  $P^T H Q$ ,  $(\Phi)_{i,i}$  is the  $i$ -th element of the diagonal matrix  $\Phi \in \mathbb{R}^{n \times n}$ , and  $(\Psi)_{j,j}$  is the  $j$ -th element of the diagonal matrix  $\Psi \in \mathbb{R}^{d \times d}$ .

As we know,  $(\Phi)_{i,i} \geq 0$  and  $(\Psi)_{j,j} \geq 0$ . In the meantime,  $\rho_1$  and  $\rho_2$  are greater than zero in practice, so the denominator in the equation above is greater than zero. After obtaining  $D$ , we can easily calculate  $X$  as

$$\begin{aligned} P^T X Q &= D \in \mathbb{R}^{n \times d} \\ P P^T X Q Q^T &= X = P D Q^T \in \mathbb{R}^{n \times d} \end{aligned} \quad (18)$$

Therefore,  $X^{(t+1)} = X = P D Q^T$ .

(2) Compute the subproblem about  $Y^{(t+1)}$ :

Further we calculate the subproblem about  $Y^{(t+1)}$ . When the other variables are fixed with the former iteration result  $(X^{(t+1)}, Z^{(t)}, \Lambda_1^{(t)}, \Lambda_2^{(t)})$ , the subproblem about  $Y^{(t+1)}$  is as

$$\begin{aligned} Y^{(t+1)} &= \operatorname{argmin}_{Y \in \mathbb{R}^{n \times d}} \mathcal{L}_{\rho_1, \rho_2}(X^{(t+1)}, Y, Z^{(t)}, \Lambda_1^{(t)}, \Lambda_2^{(t)}) \\ &= \operatorname{argmin}_{Y \in \mathbb{R}^{n \times d}} \left\{ \alpha \|Y\|_{2,1} + \frac{\rho_1}{2} \|X^{(t+1)} - Y + \frac{1}{\rho_1} \Lambda_1^{(t)}\|_F^2 \right\} \end{aligned} \quad (19)$$

In order to solve the subproblem (19), we first decouple it as

$$\begin{aligned} Y^{(t+1)} &= \operatorname{argmin}_{Y_i \in \mathbb{R}^d, i=1,2,\dots,n} \left\{ \alpha \sum_{i=1}^n \|Y_i\|_2 \right. \\ &\quad \left. + \frac{\rho_1}{2} \sum_{i=1}^n \|Y_i - (X^{(t+1)} + \frac{1}{\rho_1} \Lambda_1^{(t)})_i\|_2^2 \right\} \end{aligned} \quad (20)$$

where  $Y_i \in \mathbb{R}^d$  and  $(X^{(t+1)} + \frac{1}{\rho_1} \Lambda_1^{(t)})_i \in \mathbb{R}^d$  ( $i = 1, 2, \dots, n$ ) are the  $i$ -row of the matrices  $Y \in \mathbb{R}^{n \times d}$  and  $X^{(t+1)} + \frac{1}{\rho_1} \Lambda_1^{(t)} \in \mathbb{R}^{n \times d}$ , respectively.

Define  $Y_i^{(t+1)} \in \mathbb{R}^d$  as the  $i$ -row of the matrix  $Y^{(t+1)} \in \mathbb{R}^{n \times d}$  ( $i = 1, 2, \dots, n$ ), then, we have

$$\begin{aligned} Y_i^{(t+1)} &= \operatorname{argmin}_{Y_i \in \mathbb{R}^d} \left\{ \alpha \|Y_i\|_2 \right. \\ &\quad \left. + \frac{\rho_1}{2} \|Y_i - (X^{(t+1)} + \frac{1}{\rho_1} \Lambda_1^{(t)})_i\|_2^2 \right\} \end{aligned} \quad (21)$$

The problem (21) can be solved by the following lemma:

**Lemma 1** For any  $\lambda > 0$  and  $\mathbf{u} \in \mathbb{R}^d$ , the minimizer of

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \|\mathbf{w}\|_2 \right\} \quad (22)$$

is given by

$$\mathbf{w}^* = \begin{cases} \mathbf{0} & \text{if } \|\mathbf{u}\|_2 \leq \lambda \\ (1 - \frac{\lambda}{\|\mathbf{u}\|_2})\mathbf{u} & \text{if } \|\mathbf{u}\|_2 > \lambda \end{cases} \quad (23)$$

Based on this lemma, we can obtain the optimal  $\mathbf{Y}_i^{(t+1)} \in \mathbb{R}^d$  as

$$\mathbf{Y}_i^{(t+1)} = \begin{cases} \mathbf{0} & \text{if } \|\mathbf{S}_i^{(t)}\|_2 \leq \frac{\alpha}{\rho_1} \\ (1 - \frac{\alpha}{\rho_1 \|\mathbf{S}_i^{(t)}\|_2})\mathbf{S}_i^{(t)} & \text{if } \|\mathbf{S}_i^{(t)}\|_2 > \frac{\alpha}{\rho_1} \end{cases} \quad (24)$$

where  $\mathbf{S}_i^{(t)}$  is the  $i$ -th row of matrix  $\mathbf{S}^{(t)}$  ( $i = 1, 2, \dots, n$ ), and  $\mathbf{S}^{(t)} = \mathbf{X}^{(t+1)} + \frac{1}{\rho_1} \mathbf{\Lambda}_1^{(t)} \in \mathbb{R}^{n \times d}$ .

(3) Compute the subproblem about  $\mathbf{Z}^{(t+1)}$ :

Further we calculate the subproblem about  $\mathbf{Z}^{(t+1)}$ . When the other variables are fixed with the former iteration result  $(\mathbf{X}^{(t+1)}, \mathbf{Y}^{(t+1)}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)})$ , the subproblem about  $\mathbf{Z}^{(t+1)}$  is as

$$\begin{aligned} \mathbf{Z}^{(t+1)} &= \operatorname{argmin}_{\mathbf{Z} \in \mathbb{R}^{d \times n}} \mathcal{L}_{\rho_1, \rho_2}(\mathbf{X}^{(t+1)}, \mathbf{Y}^{(t+1)}, \mathbf{Z}, \mathbf{\Lambda}_1^{(t)}, \mathbf{\Lambda}_2^{(t)}) \\ &= \operatorname{argmin}_{\mathbf{Z} \in \mathbb{R}^{d \times n}} \left\{ \beta \|\mathbf{Z}\|_{2,1} + \frac{\rho_2}{2} \|(\mathbf{X}^{(t+1)})^T - \mathbf{Z} + \frac{1}{\rho_2} \mathbf{\Lambda}_2^{(t)}\|_F^2 \right\} \end{aligned} \quad (25)$$

Similar to solve (19), the optimal  $\mathbf{Z}^{(t+1)}$  can be easily obtained by

$$\mathbf{Z}_j^{(t+1)} = \begin{cases} \mathbf{0} & \text{if } \|\mathbf{T}_j^{(t)}\|_2 \leq \frac{\beta}{\rho_2} \\ (1 - \frac{\beta}{\rho_2 \|\mathbf{T}_j^{(t)}\|_2})\mathbf{T}_j^{(t)} & \text{if } \|\mathbf{T}_j^{(t)}\|_2 > \frac{\beta}{\rho_2} \end{cases} \quad (26)$$

where  $\mathbf{T}_j^{(t)}$  is the  $j$ -th row of matrix  $\mathbf{T}^{(t)}$  ( $j = 1, 2, \dots, d$ ), and  $\mathbf{T}^{(t)} = (\mathbf{X}^{(t+1)})^T + \frac{1}{\rho_2} \mathbf{\Lambda}_2^{(t)} \in \mathbb{R}^{d \times n}$ .

## 5. Complexity Analysis

The main computation in each iteration comes from updating  $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $\mathbf{Z}$  and dual variables  $\mathbf{\Lambda}_1$  and  $\mathbf{\Lambda}_2$ . The update steps of the dual variables refer to one matrix multiplication and several matrix additions whose computational complexity is  $\mathcal{O}(n^2d)$ . For updating  $\mathbf{X}$ , it refers to several matrix multiplications and two eigenvalue decompositions, which costs  $\mathcal{O}(n^3 + n^2d + d^2n + d^3)$ . For updating  $\mathbf{Y}$ , the complexity is of order  $\mathcal{O}(nd)$ . Updating  $\mathbf{Z}$  needs  $\mathcal{O}(n^2d)$ . Therefore, the total computational complexity in each iteration is  $\mathcal{O}(n^3 + n^2d + d^2n + d^3)$ .

## 6. Conclusion

This work provides a theoretical foundation for a unified framework that uses CUR matrix decomposition to perform active learning (selecting informative instances) and feature selection simultaneously. The analysis proves that selecting the most influential rows and columns via their optimization processing offers strong performance guarantees. The method is interpretable as it uses actual data points and features, offering a great potential for high-dimensional data analysis [7, 8, 20, 29], group fairness analysis [19, 27], and outlier or anomaly detection tasks [10, 15, 18].

## References

- [1] Hervé Abdi. Singular value decomposition (svd) and generalized singular value decomposition. *Encyclopedia of Measurement and Statistics*, 907(912):44, 2007.
- [2] Afshin Ahmadi, Felice Manganiello, Amin Khademi, and Melissa C Smith. A parallel jacobi-embedded gauss-seidel method. *IEEE Transactions on Parallel and Distributed Systems*, 32(6): 1452–1464, 2021.
- [3] Christos Boutsidis and David P Woodruff. Optimal cur matrix decompositions. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 353–362, 2014.
- [4] Davide Cacciarelli and Murat Kulahci. Active learning for data streams: a survey. *Machine Learning*, 113(1):185–239, 2024.
- [5] Quézia Cavalcante and Milton J Porsani. Low-rank seismic data reconstruction and denoising by cur matrix decompositions. *Geophysical Prospecting*, 70(2):362–376, 2022.
- [6] Zhong Chen, Zhide Fang, Wei Fan, Andrea Edwards, and Kun Zhang. Cstg: An effective framework for cost-sensitive sparse online learning. In *Proceedings of the 2017 SIAM International Conference on Data Mining*, pages 759–767. SIAM, 2017.
- [7] Zhong Chen, Zhide Fang, Victor Sheng, Jiabin Zhao, Wei Fan, Andrea Edwards, and Kun Zhang. Adaptive robust local online density estimation for streaming data. *International Journal of Machine Learning and Cybernetics*, 12(6):1803–1824, 2021.
- [8] Zhong Chen, Huixin Zhan, Victor Sheng, Andrea Edwards, and Kun Zhang. Projection dual averaging based second-order online learning. In *2022 IEEE International Conference on Data Mining*, pages 51–60. IEEE, 2022.
- [9] Zhong Chen, Huixin Zhan, Victor Sheng, Andrea Edwards, and Kun Zhang. Proximal cost-sensitive sparse group online learning. In *2022 IEEE International Conference on Big Data*, pages 495–504. IEEE, 2022.
- [10] Zhong Chen, Victor Sheng, Andrea Edwards, and Kun Zhang. An effective cost-sensitive sparse online learning framework for imbalanced streaming data classification and its application to online anomaly detection. *Knowledge and Information Systems*, 65(1):59–87, 2023.
- [11] Zhong Chen, Yi He, Di Wu, Huixin Zhan, Victor Sheng, and Kun Zhang. Robust sparse online learning for data streams with streaming features. In *Proceedings of the 2024 SIAM International Conference on Data Mining*, pages 181–189. SIAM, 2024.
- [12] Zhong Chen, Yi He, Di Wu, Liudong Zuo, Keren Li, Wenbin Zhang, and Zhiqiang Deng.  $\ell_{1,2}$ -norm and cur decomposition based sparse online active learning for data streams with streaming features. In *2024 IEEE International Conference on Big Data*, pages 384–393. IEEE, 2024.
- [13] Zhong Chen, Victor Sheng, Andrea Edwards, and Kun Zhang. Cost-sensitive sparse group online learning for imbalanced data streams. *Machine Learning*, 113(7):4407–4444, 2024.

- [14] Zhong Chen, Yi He, Di Wu, Wenbin Zhang, and Zhiqiang Deng.  $\ell_{1,\infty}$  mixed norm promoted row sparsity for fast online cur decomposition learning in varying feature spaces. In *Proceedings of the 2025 SIAM International Conference on Data Mining*, pages 124–133. SIAM, 2025.
- [15] Zhong Chen, Yi He, Di Wu, Chen Zhao, and Meikang Qiu. A novel sparse active online learning framework for fast and accurate streaming anomaly detection over data streams. In *International Joint Conferences on Artificial Intelligence*, pages 2740–2748, 2025.
- [16] Wei Deng, Ming-Jun Lai, Zhimin Peng, and Wotao Yin. Parallel multi-block admm with  $\mathcal{O}(1/k)$  convergence. *Journal of Scientific Computing*, 71(2):712–736, 2017.
- [17] Claudio Gambella and Andrea Simonetto. Multiblock admm heuristics for mixed-binary optimization on classical and quantum computers. *IEEE Transactions on Quantum Engineering*, 1:1–22, 2020.
- [18] Heng Lian, Yi He, Di Wu, Zhong Chen, Xingquan Zhu, and Xindong Wu. Online outlier detection in open feature spaces. *IEEE Transactions on Knowledge and Data Engineering*, 2025.
- [19] Heng Lian, Chen Zhao, Zhong Chen, Xingquan Zhu, My T Thai, and Yi He. Metric-agnostic continual learning for sustainable group fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 18648–18657, 2025.
- [20] Cheng Liang, Di Wu, Yi He, Teng Huang, Zhong Chen, and Xin Luo. Mma: Multi-metric-autoencoder for analyzing high-dimensional and incomplete data. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 3–19. Springer, 2023.
- [21] Sanmin Liu, Shan Xue, Jia Wu, Chuan Zhou, Jian Yang, Zhao Li, and Jie Cao. Online active learning for drifting data streams. *IEEE Transactions on Neural Networks and Learning Systems*, 34(1):186–200, 2021.
- [22] Andrzej Maćkiewicz and Waldemar Ratajczak. Principal components analysis (pca). *Computers & Geosciences*, 19(3):303–342, 1993.
- [23] Michael W Mahoney and Petros Drineas. Cur matrix decompositions for improved data analysis. *Proceedings of the National Academy of Sciences*, 106(3):697–702, 2009.
- [24] Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. Active learning by acquiring contrastive examples. *arXiv preprint arXiv:2109.03764*, 2021.
- [25] Jialei Wang, Peilin Zhao, Steven CH Hoi, and Rong Jin. Online feature selection and its applications. *IEEE Transactions on Knowledge and Data Engineering*, 26(3):698–710, 2013.
- [26] Jing Wang, Meng Wang, Peipei Li, Luoqi Liu, Zhongqiu Zhao, Xuegang Hu, and Xindong Wu. Online feature selection with group structure analysis. *IEEE Transactions on Knowledge and Data Engineering*, 27(11):3029–3041, 2015.
- [27] Zichong Wang, Jocelyn Dzuong, Xiaoyong Yuan, Zhong Chen, Yanzhao Wu, Xin Yao, and Wenbin Zhang. Individual fairness with group awareness under uncertainty. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 89–106. Springer, 2024.

- [28] Di Wu, Yi He, Xin Luo, and MengChu Zhou. A latent factor analysis-based approach to online sparse streaming feature selection. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 52(11):6744–6758, 2021.
- [29] Di Wu, Shengda Zhuo, Yu Wang, Zhong Chen, and Yi He. Online semi-supervised learning with mix-typed streaming features. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 4720–4728, 2023.
- [30] Xindong Wu, Kui Yu, Wei Ding, Hao Wang, and Xingquan Zhu. Online feature selection with streaming features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(5): 1178–1192, 2012.

## Appendix

The key steps of the CUR algorithm are summarized in Algorithm 1. We can extend our method to the kernel version by defining a new data representation to incorporate the kernel information.

---

### Algorithm 1 The CUR Decomposition Algorithm

---

**Online input:** the data matrix  $\mathbf{W} \in \mathbb{R}^{d \times n}$ , parameters,  $\alpha > 0$  and  $\beta > 0$ .

**Online output:** the matrix,  $\mathbf{X}^{(t)} \in \mathbb{R}^{n \times d}$ .

- 1: **Initialization:**  $\mathbf{Y}^{(0)} = \mathbf{0} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{Z}^{(0)} = \mathbf{0} \in \mathbb{R}^{d \times n}$ ,  $\mathbf{\Lambda}_1^{(0)} = \mathbf{0} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{\Lambda}_2^{(0)} = \mathbf{0} \in \mathbb{R}^{d \times n}$ ,  $\rho_1 > 0$ ,  $\rho_2 > 0$ ,  $\rho = 10^{10}$ ,  $\tau > 0$ ,  $\varepsilon > 0$ ,  $t = 0$ ;
  - 2: **While not converged do**
  - 3: fix the other variables and update  $\mathbf{X}^{(t+1)}$  by  $\mathbf{X}^{(t+1)} = \mathbf{X} = \mathbf{P}\mathbf{D}\mathbf{Q}^T$  in Eq. (18);
  - 4: fix the other variables and update  $\mathbf{Y}^{(t+1)}$  by Eq. (24);
  - 5: fix the other variables and update  $\mathbf{Z}^{(t+1)}$  by Eq. (25);
  - 6: update the multipliers  $\mathbf{\Lambda}_1^{(t+1)} = \mathbf{\Lambda}_1^{(t)} + \rho_1(\mathbf{X}^{(t+1)} - \mathbf{Y}^{(t+1)})$  and  $\mathbf{\Lambda}_2^{(t+1)} = \mathbf{\Lambda}_2^{(t)} + \rho_2((\mathbf{W}^{(t+1)})^T - \mathbf{Z}^{(t+1)})$ ;
  - 7: update the parameters  $\rho_1$  and  $\rho_2$  by  $\rho_1 = \min(\tau\rho_1, \rho)$  and  $\rho_2 = \min(\tau\rho_2, \rho)$ ;
  - 8:  $t \leftarrow t + 1$ ;
  - 9: check the convergence conditions (1)  $\|\mathbf{X}^{(t)} - \mathbf{Y}^{(t)}\|_\infty < \varepsilon$ ; (2)  $\|(\mathbf{X}^{(t)})^T - \mathbf{Z}^{(t)}\|_\infty < \varepsilon$ ; and (3)  $|\frac{f(\mathbf{X}^{(t)}) - f(\mathbf{X}^{(t-1)})}{f(\mathbf{X}^{(t-1)})}| < \varepsilon$  are all satisfied, where  $f(\mathbf{X}) = \|\mathbf{W} - \mathbf{W}\mathbf{X}\mathbf{W}\|_F^2 + \alpha\|\mathbf{X}\|_{2,1} + \beta\|\mathbf{X}^T\|_{2,1}$  is the objective function value of Eq. (4) at the point  $\mathbf{X}$ .
  - 10: **end**
-