

RESEARCH ARTICLE OPEN ACCESS

An Efficient Filtering Approach for Model Estimation in Sparse Regression

Yanran Wei¹  | William Myers² | Xinwei Deng¹ 

¹Department of Statistics, Virginia Tech, Blacksburg, Virginia, USA | ²Farmer School of Business, Miami University, Oxford, Ohio, USA

Correspondence: Xinwei Deng (xdeng@vt.edu)

Received: 31 March 2024 | **Revised:** 12 September 2024 | **Accepted:** 14 September 2024

Funding: This work was supported by a Research Collaboration Program between Procter & Gamble Co. and Virginia Tech.

Keywords: data reduction | data structure | filtering approach | sparse regression

ABSTRACT

As technology advances, the scale of data generated is growing exponentially, bringing huge challenges to data storage and computation. To facilitate the computational cost while maintaining model estimation accuracy, subdata selection become important. Conventional methods, such as LASSO and ridge regression, often focus on feature selection. In contrast, methods of subsampling aim at specifying data points to be extracted. However, these subsampling methods often overlook the full consideration of the role of the response variable and its relationship with predictor variables. In this work, we propose a so-called filtering approach for model estimation (FAME) method to perform subsampling in combination with feature screening. Compared with existing methods, the proposed method can result in the subdata being smaller in size both in terms of the number of features and observations, and also the computational complexity does not increase. The proposed method can be extended to situations when the predictor is binary, the response is binary, or both are binary. The performance of FAME is evaluated in both numerical studies and real data examples.

1 | Introduction

As modern technologies advance, the capability of collecting large datasets is often beyond the analysis ability of traditional statistical methods. For example, Walmart processes over 40 Petabytes of data with millions of rows per day, aiming to obtain valuable information on customer needs. Extracting valuable information from such large-scale data is challenging, and it is beyond the capability of conventional statistical methods to deal with these massive datasets directly. There is an emerging need to conduct data reduction for such large datasets and enable efficient data analysis.

Without loss of generality, this paper concerns the dataset with an $n \times 1$ response vector and an $n \times p$ predictor matrix. Here,

a large dataset can be referred to as the situation that both n and p can be large. In the literature, there are two types of common approaches for data reduction. The first common method is subsampling of data points to get a subdata with n^* data points, where $n^* \ll n$. For example, Drineas et al. [1] developed algorithms based on statistical leverage scores of matrices. Ma, Mahoney, and Yu [2] proposed two new leveraging algorithms: shrunk leveraging estimator and unweighted leveraging estimator. The algorithmic leveraging method takes at least $O(np \log n / \epsilon^2)$ ($\epsilon \in (0, 0.5]$) time, and $O(np^2)$ for appropriate parameter setting. Besides subsampling-based methods, Wang, Yang, and Stufken [3] developed information-based optimal subdata selection (IBOSS) method selecting data points based on a certain optimal design criterion. Deldossi and Tommasi [4] proposed the method “Optimal Design Based” for observation selection based on any optimality criterion. Xie, Bai, and Ma

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2024 The Author(s). Statistical Analysis and Data Mining published by Wiley Periodicals LLC.

[5] proposed data reduction methods on D-optimality improving the computational efficiency of the online analysis. Meng et al. [6] used orthogonal Latin hypercube designs for robust regression estimation even when the underlying linear model is misspecified. Singh and Stufken [7] proposed the Combining Lasso and Subdata Selection method combining multiple LASSO and IBOSS for subdata selection. More work is extended based on IBOSS, like Cheng, Wang, and Yang [8] applied subdata selection on logistic regression models and Ai et al. [9] proposed the subsampling method for quantile regression.

Considering the data with a large dimension of predictor variables, the second common method for data reduction is to get a subdata with p^* predictors where $p^* \ll p$. Various variable selection methods have been proposed, including LASSO [10], SCAD [11], and Dantzig [12], among many others. In particular, Fan and Lv [13] proposed the sure independence screening (SIS) method to reduce the dimensionality of predictor variables with fast computation. Wang [14] proposed FP-SIS for ultrahigh-dimensional variable selection by factor analysis. Zhao et al. [15] extended SIS to a so-called preconditioned profiled independence screening (PPIS) method that can perform consistent model selection for spiked populations. Note that the variable selection approaches do not focus on the reduction of data points, while the subsampling approaches often overlook the problem of the large dimensionality of predictor variables. Moreover, the subsampling approaches based on optimal design criteria often fail to take full advantage of the information on the response variables.

In this work, we proposed a so-called filtering approach for model estimation (FAME), which efficiently selects subdata for estimating regression models. The key idea of the proposed method is to consider data reduction on both predictor variables and data points with fast computation. Specifically, we adopt the sure independence screening to effectively reduce the number of predictor variables in the subdata, and then use the information-based optimal subdata selection to reduce the number of data points in the subdata. The proposed FAME approach has several advantages over other methods. First, the dataset selected by the proposed FAME method is much smaller in both size and dimension than the original data. While the model estimation accuracy based on the reduced data is comparable to that based on the original data. Second, the computation time of the proposed FAME method is in the order of $O(np)$. It implies that there is no extra computation cost compared with other methods. Third, we consider both the response variable and predictor variable in the reduction of data points, which is different from the optimal design-based subdata selection. Finally, although this paper focuses on both the response variable and predictor variable to be continuous, the proposed FAME method can be easily extended to the case when the response variable or predictor variable is discrete.

The remainder of the paper is organized as follows. In Section 2, we present the proposed FAME method for continuous data and derive the statistical properties of the proposed FAME estimators. We extend it to situations with binary predictor variables or binary response variables. The proposed FAME method is evaluated using simulation data in Section 3 and is also examined in real data in Section 4. In Section 5, we conclude the paper with discussions on further research directions.

2 | The Proposed Methodology

In this section, we will detail the proposed method. Suppose that there are p predictor variables $x_j, j = 1, \dots, p$ and the corresponding response variable y . Let $\mathbf{D} = (\mathbf{y}, \mathbf{X})$ denote the full dataset, where $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ is an n -vector continuous response variable, and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ is an n by p input matrix, where $\mathbf{x}_i \in \mathcal{R}^p$ for $i = 1, \dots, n$. We consider the linear regression model for \mathbf{y} and \mathbf{X} as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

where $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)$ is the parameter vector, and $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$ is the error vector with $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$. When $n > p$, the estimation of parameter $\boldsymbol{\beta}$ depends on the matrix inversion $(\mathbf{X}^T \mathbf{X})^{-1}$, which can be expensive when the dimensionality p is large. In the high-dimensional setting with $p > n$, the regularization methods, such as LASSO, are often used for model estimation. It is known that the computational complexity of LASSO [16] is closely related to both the sample size n and dimensionality p . Moreover, when the sample size n is large, the data could contain redundant information in rows, which may not contribute much to the inference and estimation of $\boldsymbol{\beta}$. To overcome the above challenges, we consider an efficient data reduction in terms of reducing both the rows and columns of the original data for model estimation in sparse regression. We will start with linear models with continuous inputs in Section 2.1 and extend the proposed method for generalized linear models in Section 2.3.

2.1 | The FAME Method

To efficiently conduct data reduction, we propose an FAME in the sparse regression setting. The key idea behind the proposed FAME method is to select a subset of columns and a subset of rows based on data filtering to form the reduced dataset for model estimation. Here, the data filtering aims to filter out (i.e., remove) relatively unimportant columns or rows such that the reduced dataset contains key information with a smaller size than the original dataset.

Specifically, we first conduct a feature screening to select important columns. Denote $\mathbf{x}^{(j)}$ to be the j th column of input matrix \mathbf{X} and recall that the response vector is \mathbf{y} . Without loss of generality, we assume each column of \mathbf{X} is standardized. Motivated by the sure independent screening in Fan and Lv [13], we consider to filter out features based on each individual column's index score as $w_j = g(\mathbf{x}^{(j)}, \mathbf{y})$. Thus, we can keep a subset of columns with the highest scores. Here, the score index measures the importance of feature $\mathbf{x}^{(j)}$ by the degree of correlation between $\mathbf{x}^{(j)}$ and the response variable y for $j = 1, \dots, p$. For example, the index score w_j can be the absolute value of the regression coefficient derived from the marginal regression of \mathbf{y} on $\mathbf{x}^{(j)}$ as shown in Fan and Lv [13]. Then, we can define a subset of predictors as

$$\mathbf{M}_S = \{1 \leq q \leq h \leq p : w_q \text{ is among the } h \text{ largest values in } \mathbf{w} = (w_1, \dots, w_p)^T\} \quad (2)$$

Correspondingly, the reduced input matrix consists of $\mathbf{X}_S = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(q)}, \dots, \mathbf{x}^{(h)})$, where $\mathbf{x}^{(q)}$ is the column whose score

index belongs to M_S . Then, the complete data D are reduced to a subset with h columns as $D_S = (X_S, y)$.

Although the reduction in the number of predictors will address the computational challenges of model estimation to some extent, there may be redundant information among data points for the big dataset with large size of rows. Thus, in the second step, we filter out relatively unimportant rows in the dataset D_S . Specifically, the proposed FAME method selects data points with the aim of maximizing the Fisher information matrix based on the model in Equation (1). It is known that the Fisher information matrix based on the model in Equation (1) for the reduced dataset D_S can be written as $I(\delta) = \frac{1}{\sigma^2} \sum_{i=1}^n \delta_i x_{is} x_{is}^T$, where $x_{is} \in \mathcal{R}^h$ is the data point in D_S and $\delta_i = 1$ means the i th row in D_S is selected, otherwise not. Therefore, one could consider to select a subset of rows with size k based on the following criteria:

$$\max_{\delta=(\delta_1, \dots, \delta_n), \delta_i \in \{0,1\}} \left| \sum_{i=1}^n \delta_i x_{is} x_{is}^T \right|, \quad \text{s.t.} \quad \sum_{i=1}^n \delta_i = k \quad (3)$$

which is also known as the D-optimality criterion in the design of experiment literature Kiefer [17]. However, the optimization in Equation (3) is nontrivial due to the discrete nature of the optimization. To address this issue, we consider an upper bound of the determinant of $I(\delta)$ such that the selection of rows becomes computationally efficient. Motivated by Wang, Yang, and Stufken [3], we consider one reasonable upper bound for $|I(\delta)|$ to be,

$$|I(\delta)| \leq C \frac{k^{h+1} l^{2h}}{4^h \sigma^{2(h+1)}} \prod_{q=1}^h (\bar{x}_{(u)q} - \bar{x}_{(v)q})^2, \quad q = 1, \dots, h \quad (4)$$

where $\bar{x}_{(u)q} = \frac{1}{l} (x_{(n-l+1)q} + \dots + x_{(n)q})$ is the average of the l largest order statistic in $x^{(q)}$. Here, $x_{(1)q} \leq x_{(2)q} \leq \dots \leq x_{(n)q}$ is the order statistic of the q th column $x^{(q)}$. Similarly, $\bar{x}_{(v)q} = \frac{1}{l} (x_{(1)q} + \dots + x_{(l)q})$ is the average of the l smallest order statistic in $x^{(q)}$. Note that the original upper bound in Wang, Yang, and Stufken [3] is a special case of the upper bound in Equation (4). By considering the average of the upper-order statistics and the average of the lower-order statistic of $x^{(q)}$, it would help the selection of data points robust against the potential outliers and uncertainty on the extreme observations. Thus, we consider the proposed FAME method of selecting a subset of rows as

$$\max_{\delta=(\delta_1, \dots, \delta_n), \delta_i \in \{0,1\}} \sum_{i=1}^h \delta_i \log [\bar{x}_{(u)q} - \bar{x}_{(v)q}], \quad \text{s.t.} \quad \sum_{i=1}^n \delta_i = k \quad (5)$$

It is seen that the above optimization can be decomposed into the selection of rows based on each column, which is computationally fast. For practical implementation, we sort the predictor values on the q th column $x^{(q)}$ and select the $2l$ data points having the l largest and the l smallest order statistics. Repeat such a procedure for each column in D_S , we compose the selected subset of data as $\tilde{D}_S = (\tilde{X}_S, \tilde{y})$. Note that the reduced input matrix \tilde{X}_S is a $k \times h$ matrix. The choice of l and h is closely related to the size of the filtered data, $k = 2lh$, which will be discussed at the end of this subsection.

With the filtered data \tilde{D}_S , we consider the model estimation by using the LASSO method [10] as

$$\tilde{\beta} = \arg \min_{\beta} (\tilde{y} - \tilde{X}_S \beta)^T (\tilde{y} - \tilde{X}_S \beta) + \lambda \|\beta\|_1 \quad (6)$$

where $\lambda \geq 0$ is a tuning parameter and $\|\beta\|_1 = \sum_{j=1}^h |\beta_j|$ is the l_1 penalty. The optimal value of λ is chosen using cross-validation. Algorithm 1 summarizes the proposed FAME algorithm for sparse regression as follows.

ALGORITHM 1 | (FAME for sparse regression).

Step 1: For $1 \leq j \leq p$, calculate the score index $w_j = g(x^{(j)}, y)$;

Step 2: Form $D_S = (X_S, y)$, where X_S includes the h columns with the h largest w_j values;

Step 3: For $q = 1, \dots, h$, based on column $x^{(q)}$ in X_S , form $\tilde{D}_S = (\tilde{X}_S, \tilde{y})$ by including $2l$ data points with the l smallest values and l largest values but excluding data points already been selected;

Step 4: Obtain parameter estimation from Equation (6) using the data \tilde{D}_S .

One can see that the computation complexity of Algorithm 1 is on the order of $O(np)$, which can be much faster than the Lasso using the original dataset. There are several remarks on the specification of Algorithm 1. First, the index score w_j measures the relationship between y and $x^{(j)}$. It can be the Pearson correlation, partial correlation, or other marginal statistics, like t -statistics and p value. Based on the empirical study, we recommend using t -statistic providing robust results from Algorithm 1 compared with Pearson correlation and p value. Details will be discussed in the experimental section.

Here, we would like to remark that the choice of k , the size of reduced data, is often constrained by the practical consideration. Considering the selected columns have the most strong score index explaining the response variables, and the rule of thumb based on our empirical study (see a numerical study in Section 3) is to select 10–30 points from each column.

For the choice of h , the number of important features to be retained in the filtered data, one may follow the suggestion given by Fan and Lv [13] with $h = n/\log(n)$ when $n < p$ and $h = p/\log(p)$ when $n > p$. However, when both n and p are large, such a choice of h can also be large. Alternatively, one can determine h by using the change-point detection method based on the score index values $w_j = g(x^{(j)}, y)$, $j = 1, \dots, p$. The change-point detection method is to find the first point when the distribution of the descending-ordered score index w_j changes. Here, we adopt the pruned exact linear time (PELT) change-point detection algorithm in Killick, Fearnhead, and Eckley [18]. Specifically, suppose that there are a sorted score index $(w^{(1)}, \dots, w^{(p)})$ and τ change-points splitting the data into $\tau + 1$ segments. Denote the τ change-point locations as $L = (l_1, \dots, l_\tau)$. Then, the j th segment contains all score index between location $l_{j-1} + 1$ and l_j , denoted as $w^{(l_{j-1}+1):l_j}$. A change point is detected if we can detect a point location

j that satisfies,

$$C(\mathbf{w}^{(t+1):j}) + C(\mathbf{w}^{(j+1):T}) + K \leq C(\mathbf{w}^{(t+1):T}),$$

for all $t < j < T$ (7)

where $C(\cdot)$ is a cost function for a segment, and K is a constant. The settings of $C(\cdot)$ and K follow the algorithm in Killick, Fearnhead, and Eckley [18] and the corresponding R package *changeoint*.

2.2 | Theoretical Properties

In this section, we investigate the theoretical properties of the FAME algorithm. These properties evaluate our proposed method and give bounds on variances of estimators as well as asymptotic results.

Theorem 1. Let $\tilde{\mathbf{D}}_S$ be the subdata selected by the FAME algorithm including h columns and k data entries in total. Denote $\lambda_{\min}(\mathbf{R})$ as the smallest eigenvalues of the sample correlation matrix \mathbf{R} of $\tilde{\mathbf{D}}_S$. The determinant $|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S|$ satisfies

$$\frac{|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S|}{\frac{k^{h+1}}{4^h} \prod_{q=1}^h (\bar{x}_{(u)q} - \bar{x}_{(v)q})^2} \geq \frac{\lambda_{\min}^h(\mathbf{R})}{h^h} \times \prod_{q=1}^h \left(\frac{x_{(n-l+1)q} - x_{(l)q}}{\bar{x}_{(u)q} - \bar{x}_{(v)q}} \right)^2$$

(8)

where $x_{(n-l+1)q}$ is the l st largest order statistic and $x_{(l)q}$ is the l st smallest order statistics in $\mathbf{x}^{(q)}$. $\bar{x}_{(u)q}$ is the average of the top l order statistic, and $\bar{x}_{(v)q}$ is the l smallest order statistic.

When $\lim_{n \rightarrow \infty} \lambda_{\min}(\mathbf{R}) > 0$, then under suitable assumptions, the lower bound of Equation (8) will not converge to 0 as $n \rightarrow \infty$. $|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S|$ is of the same order as the upper bound for $|I(\delta)|$ in Equation (4), although the upper bound is hard to achieve.

Following Tibshirani [10], one can have approximations of the lasso estimator from Equation (6) as

$$\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S \approx \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1} \tilde{\mathbf{X}}_S^T \tilde{\mathbf{y}}$$

with

$$\{\Psi(\tilde{\beta}(\lambda))\} = \text{diag}(\psi_1, \psi_2, \dots, \psi_h),$$

where $\psi_j = \begin{cases} \frac{1}{|\tilde{\beta}_j(\lambda)|} & \text{if } \tilde{\beta}_j(\lambda) \neq 0 \\ 0 & \text{otherwise} \end{cases}, \text{ for } j = 1, \dots, h.$

Then, we can also establish the properties on the variance of the approximated $\tilde{\beta}$ under the proposed FAME framework.

Theorem 2. Use the same notation as in Theorem 1, when $\lambda_{\min}(\mathbf{R}) > 0$ and $\tilde{\mathbf{D}}_S$ is the reduced dataset, then the variance of the approximated $\tilde{\beta}$ of the lasso regression satisfies

$$\text{Var}(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) \leq \frac{4h\sigma^2}{k\lambda_{\min}(\mathbf{R})(\bar{x}_{(u)q} - \bar{x}_{(v)q})^2}, \text{ for } q = 1, \dots, h$$

(9)

Theorem 2 is a finite sample property of $\tilde{\beta}$ from the subdata selected via the FAME algorithm. In the cases when $\lambda = 0$ and

all variables are selected,

$$V(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) \geq \frac{4\sigma^2}{k\lambda_{\max}(\mathbf{R})(x_{nq} - x_{1q})^2}$$

(10)

where x_{nq} is the largest order statistic and x_{1q} is the smallest order statistics in $\mathbf{x}^{(q)}$. Assume that $\lim_{n \rightarrow \infty} \lambda_{\min}^h(\mathbf{R}) > 0$. As $n \rightarrow \infty$, $\tilde{\beta}$ obtained from $\tilde{\mathbf{D}}_S$ has the following property:

$$\text{Var}(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) = O_p\left(\frac{h}{k(\bar{x}_{(u)q} - \bar{x}_{(v)q})^2}\right), q = 1, \dots, h$$

This theorem holds for any values of n , r , and h and can also be used to obtain results when one or more of these values go to infinity.

Theorem 3. Denote \mathbf{D} as the full data and the estimator of linear model using \mathbf{D} as $\hat{\beta}$. Denote the upper and lower bound of $\text{Var}(\tilde{\beta}_q|\tilde{\mathbf{D}}_S)$ as $U_{\tilde{\beta}_q|\tilde{\mathbf{D}}_S}$ and $L_{\tilde{\beta}_q|\tilde{\mathbf{D}}_S}$. Similarly, the upper and lower bound of $\text{Var}(\hat{\beta}_q|\mathbf{D})$ are $U_{\hat{\beta}_q|\mathbf{D}}$ and $L_{\hat{\beta}_q|\mathbf{D}}$. Then

$$U_{\tilde{\beta}_q|\tilde{\mathbf{D}}_S} \leq U_{\hat{\beta}_q|\mathbf{D}}$$

(11)

$$L_{\tilde{\beta}_q|\tilde{\mathbf{D}}_S} \geq L_{\hat{\beta}_q|\mathbf{D}}$$

(12)

Theorem 3 shows that the estimator based on $\tilde{\mathbf{D}}_S$ can have a smaller variance range compared to the estimator based on the full data \mathbf{D} .

2.3 | Extension of Binary Data

We can also extend the proposed FAME method for noncontinuous responses or noncontinuous predictors. Here, we focus on binary response or binary predictors with two levels.

When the response variable \mathbf{y} is binary as $y_i \in \{-1, 1\}$ for data point i and the input matrix is continuous, we can use the main effect in the experimental design to construct the score index w_j . That is, the score index w_j can be calculated as the difference of the average input between two levels of the response variable.

$$w_j = \bar{\mathbf{x}}^{(j)}(y_i+) - \bar{\mathbf{x}}^{(j)}(y_i-) = \frac{1}{n^{(+)}} \sum_{y_i=1} x_{ij} - \frac{1}{n^{(-)}} \sum_{y_i=-1} x_{ij},$$

$1 \leq j \leq p, 1 \leq i \leq n;$ (13)

where $\bar{\mathbf{x}}^{(j)}(y_i+)$ is the average of $\mathbf{x}^{(j)}$ values observed at the positive response level and $\bar{\mathbf{x}}^{(j)}(y_i-)$ is the average of $\mathbf{x}^{(j)}$ values observed at the negative response level. Here, $n^{(+)}$ and $n^{(-)}$ are the number of observations with response label $y = 1$ and $y = -1$, respectively.

Since the response variable \mathbf{y} is binary, we consider the logistic regression for model estimation. Denote $\mathbf{x}_{m\bar{s}} \in \mathcal{R}^h$ is the data point and $y_m \in \{-1, 1\}$ is the response variable in $\tilde{\mathbf{D}}_S$, where $m = 1, \dots, 2l$. Then, $\Pr(\hat{\mathbf{y}}_m = 1|\mathbf{x}_{m\bar{s}}) = \frac{\exp(\mathbf{x}_{m\bar{s}}^T \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_{m\bar{s}}^T \boldsymbol{\beta})}$, where $\boldsymbol{\beta}$ is an $h \times 1$ vector of regression coefficients. A penalized likelihood

estimation can be used for parameter estimation of β by minimizing the negative log-likelihood function.

$$\hat{\beta} = \arg \min_{\beta} \{-\ell(\beta) + \lambda \|\beta\|_1\} \quad (14)$$

where $\ell(\beta) = \frac{1}{2l} \sum_{m=1}^{2l} \log [1 + \exp(-y_m \mathbf{x}_{ms}^T \beta)]$. Here, $\lambda \geq 0$ is a tuning parameter. Algorithm 2 summarizes the proposed FAME algorithm for the binary response variable.

ALGORITHM 2 | (FAME for binary response).

-
- Step 1:** For $1 \leq j \leq p$, calculate the score index $w_j = g(\mathbf{x}^{(j)}, \mathbf{y})$ in Equation (13);
- Step 2:** Form $D_S = (\mathbf{X}_S, \mathbf{y})$, where \mathbf{X}_S includes h columns with the h largest w_j values;
- Step 3:** For $p = 1, \dots, h$, based on column $x^{(p)}$ in \mathbf{X}_S , form $\tilde{D}_S = (\tilde{\mathbf{X}}_S, \tilde{\mathbf{y}})$ by including $2l$ data points with the l smallest values and l largest values but excluding data points already been selected;
- Step 4:** Obtain model estimation from Equation (14).
-

The second extension is to consider the FAME method for data with binary predictor variables with two levels as “+” and “-.” For the selection of columns, we then can define the score index based on the difference in the average response between two levels of the predictor variable. That is,

$$w_j = \bar{y}_i(\mathbf{x}^{(j)+}) - \bar{y}_i(\mathbf{x}^{(j)-}) = \frac{1}{n^{(+)}} \sum_{x_{ij}=1} y_i - \frac{1}{n^{(-)}} \sum_{x_{ij}=-1} y_i, \quad (15)$$

$$1 \leq j \leq p, 1 \leq i \leq n$$

where $\bar{y}_i(\mathbf{x}^{(j)+})$ is the average of y_i values observed at the positive predictor level and $\bar{y}_i(\mathbf{x}^{(j)-})$ is the average of y_i values observed at the negative predictor level for the j th covariate. Here, $n^{(+)}$ and $n^{(-)}$ are the number of observations with the predictor $\mathbf{x}^{(j)} = 1$ and $\mathbf{x}^{(j)} = -1$, respectively.

For the selection of rows, however, the use of criterion in Equation (4) become improper since

$$|I(\delta)| \leq C \frac{k^{h+1} l^{2h}}{4^h \sigma^{2(h+1)}} \quad (16)$$

where C is a positive constant. Clearly, the order statistics of each column in the input matrix do not work for selecting subdata to maximize the information matrix. To overcome this issue, for each filtered column $\mathbf{x}^{(q)}$, sorting predictor values and denote the location of the value closest to $\bar{y}_i(\mathbf{x}^{(q)+})$ as $\gamma^{(q)}(+)$ and the location of the value closest to $\bar{y}_i(\mathbf{x}^{(q)-})$ as $\gamma^{(q)}(-)$, denoting $\lfloor \cdot \rfloor$ as the floor division operator. When l is odd, selecting the l data points with response variable closest to $\bar{y}_i(\mathbf{x}^{(q)+})$ denoted as $y_{\gamma^{(q)}(+)-\lfloor \frac{l}{2} \rfloor; \gamma^{(q)}(+)+\lfloor \frac{l}{2} \rfloor}(\mathbf{x}_q^+)$ and l data points with response variable closest to $\bar{y}_i(\mathbf{x}^{(q)-})$ denoted as $y_{\gamma^{(q)}(-)-\lfloor \frac{l}{2} \rfloor; \gamma^{(q)}(-)+\lfloor \frac{l}{2} \rfloor}(\mathbf{x}_q^-)$. When l is even, selecting the l data points $y_{\gamma^{(q)}(+)-\lfloor \frac{l}{2} \rfloor+1; \gamma^{(q)}(+)+\lfloor \frac{l}{2} \rfloor}(\mathbf{x}_q^+)$ and l data points $y_{\gamma^{(q)}(-)-\lfloor \frac{l}{2} \rfloor+1; \gamma^{(q)}(-)+\lfloor \frac{l}{2} \rfloor}(\mathbf{x}_q^-)$.

Algorithm 3 summarizes the proposed FAME algorithm for binary predictor variable. It is easy to see that the

computing time is the same as in Algorithm 1, which is $O(np)$.

ALGORITHM 3 | (FAME for binary predictor variables).

-
- Step 1:** For $1 \leq j \leq p$, calculate the score index $w_j = g(\mathbf{x}^{(j)}, \mathbf{y})$ in Equation (15);
- Step 2:** Form $D_S = (\mathbf{X}_S, \mathbf{y})$, where \mathbf{X}_S includes the h columns with the h largest w_j values;
- Step 3:** For $q = 1, \dots, h$, based on column $x^{(q)}$ in \mathbf{X}_S , form $\tilde{D}_S = (\tilde{\mathbf{X}}_S, \tilde{\mathbf{y}})$ by including $2l$ data points with the l values around mean of y_i when x_{iq} is positive and l values around mean of y_i when x_{iq} is negative but excluding data points already been selected;
- Step 4:** Obtain coefficient estimation from Equation (6).
-

Corollary 1. For the subdata selected by the FAME algorithm with binary predictor variables, denote \mathbf{R} as the sample correlation matrix of \tilde{D}_S , the determinant $|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S|$ satisfies

$$|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S| \geq Ck \left(\frac{k \lambda_{\min}(\mathbf{R})}{h} \right)^h \quad (17)$$

where C is a constant.

The lower bound of the information matrix of subdata \tilde{D}_S is determined by the number of columns selected in the subdata as well as the smallest eigenvalue of the correlation matrix of the subdata. Based on Theorem 3, we can know that this bound is larger than the lower bound of the information matrix of the full data and will not converge to 0 as $n \rightarrow \infty$ when $\lim_{n \rightarrow \infty} \lambda_{\min}(\mathbf{R}) > 0$.

Corollary 2. When $\lambda_{\min}(\mathbf{R}) \geq 0$, for the reduced dataset \tilde{D}_S with binary predictors,

$$\text{Var}(\tilde{\beta}_q | \tilde{D}_S) \leq \frac{Ch\sigma^2}{k\lambda_{\min}(\mathbf{R})}, \text{ for } q = 1, \dots, h \quad (18)$$

The upper bound of $\text{Var}(\tilde{\beta}_h)$ is not related to the binary predictors as in the situation of continuous predictors.

Corollary 3. As $n \rightarrow \infty$, assuming $\lim_{n \rightarrow \infty} \lambda_{\min}(\mathbf{R}) > 0$,

$$\text{Var}(\tilde{\beta}_q | \tilde{D}_S) = O_p \left(\frac{h}{k} \right), q = 1, \dots, j \quad (19)$$

When h and k are fixed, it is seen that $\text{Var}(\tilde{\beta}_q | \tilde{D}_S) = O_p(1)$.

3 | Simulation

In this section, we will evaluate the performance of the proposed FAME method using simulated data with continuous predictors or discrete predictors. To estimate the coefficient parameter $\tilde{\beta}$, we compare the proposed FAME method with five benchmark methods: IBOSS, IBOSS_LASSO, SIS, SIS_LASSO, and LASSO.

1. FAME: The proposed FAME algorithm;
2. IBOSS: Benchmark method selecting predictors randomly and selecting features using the IBOSS method [18]. The subdata are fitted with multivariate linear regression;

3. IBOSS_LASSO: Benchmark method selecting predictors randomly and selecting features using IBOSS method [18]. The subdata are fitted with regularized regression;
4. SIS: Benchmark method selecting predictors using the SIS method [13] with all data points. The subdata are fitted with multivariate linear regression;
5. SIS_LASSO: Benchmark method selecting predictors using the SIS method [13] with all data points. The subdata are fitted with regularized regression;
6. LASSO: Benchmark method on the original data fitted with regularized regression.

The subdata selected by each method are shown in Figure 1. The light color in the plot represents selected columns or rows, while dark color represents selected subdata.

Subdata are split into 70% training data and 30% testing data. Parameters are estimated on training data and the performance of methods is evaluated by empirical mean squared errors (MSEs) and signal-to-noise ratio (SNR) on testing data. $MSE = \frac{1}{n_t S} \sum_{s=1}^S \|\mathbf{y}^{(s)} - \hat{\mathbf{y}}^{(s)}\|^2$, where S is the number of simulation and n_t is the number of data points in the test dataset. $\mathbf{y}^{(s)}$ is the true response value, and $\hat{\mathbf{y}}^{(s)}$ is the predicted response value in each simulation. $SNR = \frac{1}{n} \sum_{s=1}^S \frac{\text{var}(\hat{\mathbf{y}}^{(s)})}{\hat{\sigma}^{(s)}}$, where $\hat{\sigma}^{(s)}$ is the variance of the predicted residuals.

3.1 | Cases With Continuous Predictors

Data are generated based on Equation (1) with 2000 data points, 5 true predictors among 1000 features and $\sigma^2 = 16$. Covariates matrix \mathbf{X} is generated from multivariate normal distribution $\mathbf{X} \sim N(\mathbf{0}, \Sigma)$, where Σ is a covariance matrix with $\Sigma_{ij} = 0.5^{|i-j|}$ under dependent situation and $\Sigma = \mathbf{I}$ under independent situation. $\mathbb{I}(\cdot)$ is the indicator function. The simulation is repeated $S = 50$ times, and each column is standardized. The selected subdata has 200 data points with five columns or 25 columns. Besides MSE, methods are also compared in terms of the standard deviation of MSE, true positive rate (TPR), true negative rate (TNR), and computational time (TIME), where $TPR = \text{true positive}/(\text{true positive} + \text{false negative})$ and $TNR = \text{true negative}/(\text{true negative} + \text{false positive})$.

Table 1 summarizes the estimation results of selecting 200 data points from five columns. For both independent and dependent covariates cases, the proposed FAME method beats the benchmark covariates methods with smaller MSE. Compared with IBOSS with MSE, the proposed methods selected the most important features instead of randomly selecting, which contributes to decreasing the prediction MSE from 19.49 to 16.90. Compared with SIS, the proposed FAME method uses t-statistic as the score index instead of componentwise regression. The MSE of SIS_LASSO is 37.10 which is much larger than the MSE of FAME. This big difference in prediction accuracy shows that the proposed FAME method has higher accuracy in selecting the most informative features

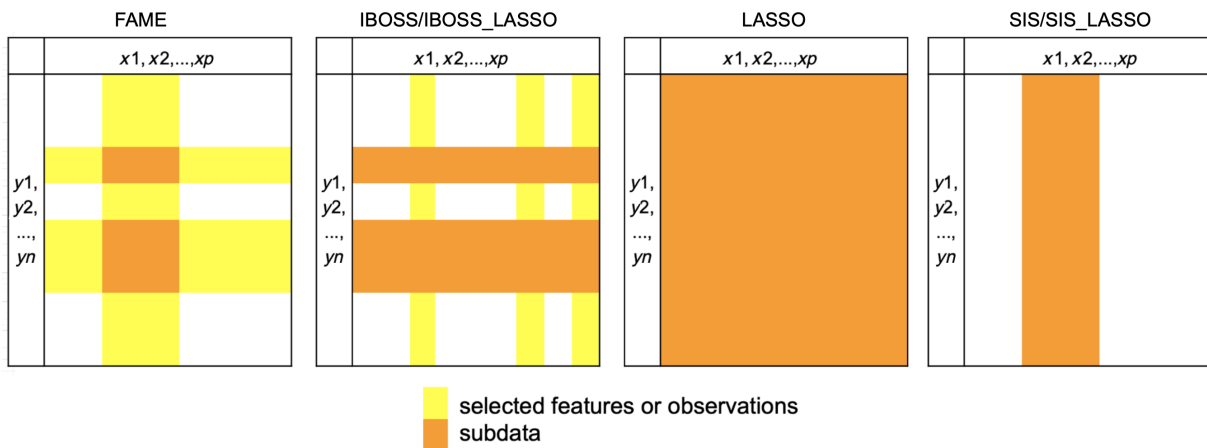


FIGURE 1 | An illustration of the proposed method and three benchmark methods.

TABLE 1 | Results of filtering 200 data points from five columns.

		FAME	IBOSS	IBOSS_LASSO	SIS	SIS_LASSO	LASSO
Independent covariates	MSE	16.90	1.35e + 07	19.49	37.10	37.10	16.43
	TPR	0.86	1.00	0.90	0.64	0.63	0.96
	TNR	1.00	0.80	0.97	1.00	1.00	0.98
	TIME	0.05	0.05	0.26	0.002	0.06	2.49
Dependent covariates	MSE	18.19	6.64e + 04	18.85	44.81	44.85	16.34
	TPR	0.82	1.00	0.90	0.63	0.61	0.96
	TNR	1.00	0.80	0.97	1.00	1.00	0.97
	TIME	0.05	0.05	0.26	0.002	0.06	1.79

compared with SIS. LASSO used all the data for prediction and has the smallest MSE of all six methods. However, the MSE of FAME is similar to that of LASSO with less number of both predictors and data points.

IBOSS, IBOSS_LASSO, and LASSO used all features for parameter estimation while FAME only utilized five features. The TPR of FAME is smaller and TNR is higher than those three methods. SIS already selected several features before regularized regression, which reduces the number of features, so the MSE of SIS_LASSO is similar to that of SIS which is 37.10. However, the regularized regression improves the IBOSS methods by reducing the MSE from $1.35e + 07$ to 19.49, leading to a decrease in TPR from 1 to 0.9.

The proposed FAME method generates a subdata with the smallest size by reducing both the number of features and data points. The computation time is also small compared to LASSO using all data and is similar to other benchmark methods. SIS has the smallest computation time because of the small number of features and no parameter tuning for regularized regression.

The result is similar when selecting 25 features as shown in Table 2. The proposed FAME method has MSE similar to that of LASSO and smaller than other benchmark methods. The prediction accuracy is high, the size of the subdata generated by FAME is the smallest among all methods and the computation costs is also low. For the simulated dataset, only five features out of 1000 are true covariates, so compared with Table 1, the MSE of FAME

does not improve, while SIS has a smaller MSE because more important features are selected.

Table 3 summarizes the variance of the estimation of selecting 200 data points from 25 or five columns. When 25 columns are selected from 1000 independent predictors, the predicted variance of all methods is smaller than that of the true variance besides the IBOSS method. The proposed FAME method has the largest SNR besides the LASSO method. The above phenomenon can also be found when the predictors are dependent or when five columns are selected from 1000 predictors. The proposed FAME method has a prediction variance close to the true variance and an SNR close to that of LASSO. The subdata generated by FAME are the smallest, and the computation time is small.

To determine the optimal number of columns in the subdata, we apply the change-point detection method described in Section 2. As shown in Figure 2, the first location of the change point detected is on the fifth column. One would recommend the number of features included in the subdata to be five, which is the number of true coefficients when generating the original data. The performance of the proposed FAME algorithm does not seem to improve much when more than five columns are included in the subdata.

Note that the number of data points contained in the subdata is 200 in this simulation study, where 40 points from each column are selected. We also examine the performance of the prediction MSE concerning the number of points in the subdata in Figure 3. As shown in Figure 3, for the dependent case, the

TABLE 2 | Results of filtering 200 data points from 25 columns.

		FAME	IBOSS	IBOSS_LASSO	SIS	SIS_LASSO	LASSO
Independent covariates	MSE	17.19	1.60e + 08	19.52	28.38	28.15	16.43
	TPR	0.90	1.00	0.90	0.72	0.72	0.96
	TNR	1.00	0.80	0.97	0.98	1.00	0.98
	TIME	0.06	0.05	0.26	0.004	0.07	2.44
Dependent covariates	MSE	17.63	1.05e + 06	18.75	28.91	28.83	16.35
	TPR	0.87	1.00	0.91	0.73	0.73	0.96
	TNR	0.99	0.80	0.97	0.98	0.99	0.97
	TIME	0.05	0.05	0.23	0.004	0.07	1.78

TABLE 3 | Results of filtering simulated subdata on prediction variance.

Method	25/1000						5/1000					
	Independent			Dependent			Independent			Dependent		
	Prediction variance	True variance	SNR	Prediction variance	True variance	SNR	Prediction variance	True variance	SNR	Prediction variance	True variance	SNR
FAME	5.40		12.94	6.95		59.70	5.45		13.78	6.96		56.58
IBOSS	18.88		1.05	13.85		1.35	16.42		1.07	11.11		1.27
IBOSS_LASSO	5.30	5.53	10.35	6.94	6.99	55.69	5.30	5.53	10.38	6.94	6.99	55.26
SIS	5.41		10.32	6.96		45.98	5.37		8.74	6.94		35.49
SIS_LASSO	5.38		10.14	6.94		45.51	5.35		8.62	6.93		34.94
LASSO	5.40		13.58	6.96		64.75	5.40		13.58	6.96		64.75

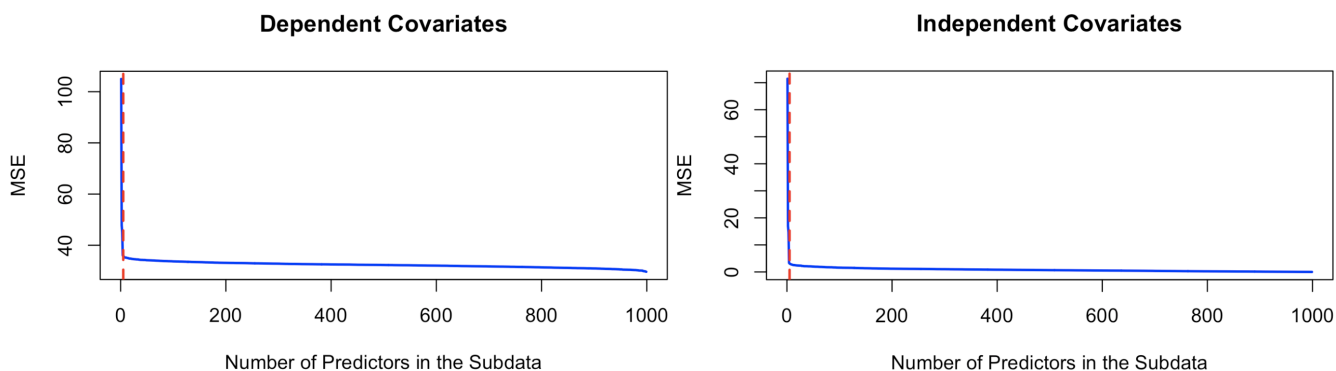


FIGURE 2 | Change-point detection using t-statistics for the simulated dataset.

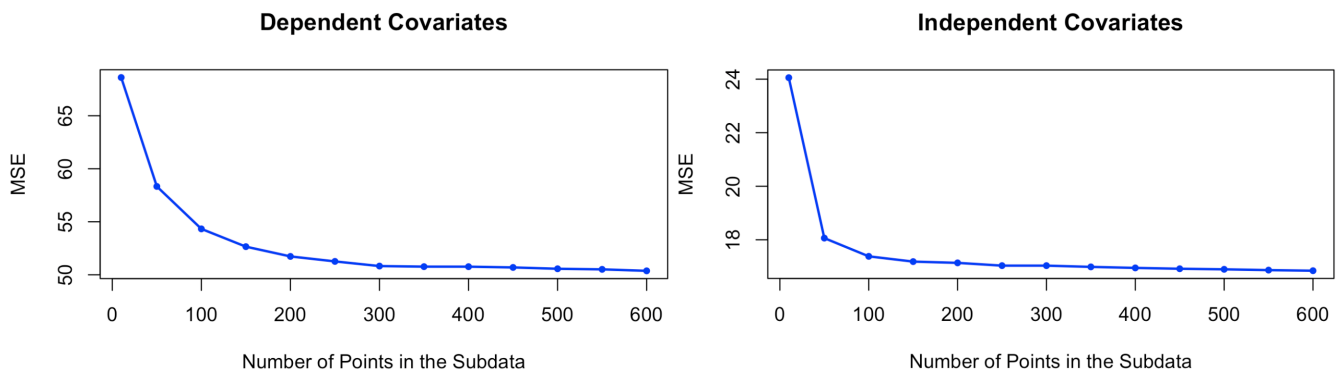


FIGURE 3 | MSE and number of points in the subdata.

MSE drops significantly when the number of data points in the subdata increases from 10 to 200. For the independent case, the turning point of MSE occurs when the subdata have around 100 data points. It is seen that the MSE does not improve much when we continuously increase the size of the subdata. Without loss of generality, one would recommend that the rule of thumb is to choose at least 10 points from each selected column to form the subdata.

3.2 | Cases With Binary Predictors

The simulated dataset has 700 observations and 700 predictors, of which five are true predictors generated from the normal distribution. The predictors came from a binomial distribution with a probability between a range of 0.3 to 0.7, representing that the covariates can be imbalanced. The response variable is generated based on Equation (1), where ϵ is the distributed normally with a mean of 0 and a standard deviation of 0.3. The simulation is repeated 50 times, and empirical mean squared error (MSE) is the measurement metric.

The simulation result is shown in Table 4. The proposed method has the smallest MSE compared with all benchmark methods, besides LASSO using all data. It shows the superior performance of feature selection and data points selection in FAME. Compared with IBOSS and IBOSS_Lasso, FAME selected the most important features and data points. Compared with SIS and SIS_Lasso, the main effect is more effective in selecting the most important features. The TPR of SIS and SIS_LASSO is only around 0.16,

while that for FAME is around 0.8. IBOSS and LASSO used all the data, so the TPR of these two methods is almost 100% accurate. The subdata generated by FAME are the smallest among all methods with low computation cost.

4 | Case Studies

4.1 | ARCENE Data in Cancer Study

The ARCENE dataset is obtained from the UCI Machine Learning Repository, which was published by the National Cancer Institute (NCI) and the Eastern Virginia Medical School (EVMS). The objective of this study is to distinguish cancer based on patterns which are continuous input variables. It is a classification problem with binary output. Among the 10,000 predictors, 7000 are real variables while the other 3000 variables are random probes. Both the training and testing dataset has 100 observations. The estimation result is measured by classification error f/n , where f is the number of sample cases incorrectly classified, and n is the total number of sample cases.

Table 5 summarizes the estimation results of the proposed and benchmark methods. When the subdata has five columns, the classification error of the proposed FAME method is the smallest among all methods. Compared with IBOSS and IBOSS_LASSO, the proposed method selected the most several important features. Compared with SIS and SIS_LASSO, the main effect performs well in ranking feature importance. Although LASSO used all features, its MSE is larger than that of FAME.

TABLE 4 | Estimating result for selecting five columns and 50 data points for dataset with binary predictors.

	FAME	IBOSS	IBOSS_LASSO	SIS	SIS_LASSO	LASSO
MSE	0.23	199.73	0.24	0.97	0.96	0.10
TPR	0.78	1.00	0.76	0.17	0.16	0.96
TNR	1.00	0.94	0.97	1.00	1.00	0.96
TIME	0.04	0.01	0.12	0.001	0.05	0.33

TABLE 5 | Results of filtering subdata with 50 observations from five or 10 columns of ARCENE data.

	c = 5		c = 10	
	Classification error	Time	Classification error	Time
FAME	0.26	0.08	0.28	0.25
IBOSS	0.47	3.88	0.48	3.61
IBOSS_LASSO	0.36	0.82	0.37	0.74
SIS	0.35	0.002	0.33	0.004
SIS_LASSO	0.42	0.84	0.40	1.85
LASSO	0.33	2.81	0.33	2.57

TABLE 6 | Result of filtering subdata with 50 observations from five columns of P&G data.

Method	R068_1min				R068_5min			
	5/59		5/3009		5/59		5/3009	
	MSE (SD)	Time (SD)	MSE (SD)	Time (SD)	MSE (SD)	Time (SD)	MSE (SD)	Time (SD)
FAME	1.87 (0.55)	0.05 (0.01)	2.05 (0.42)	0.04 (0.01)	2.94 (0.60)	0.05 (0.01)	3.22 (0.81)	0.05 (0.003)
IBOSS	2.39e + 07 (1.69e + 08)	0.002 (0.002)	4.04e + 08 (1.97e + 09)	0.03 (0.005)	1.24e + 05 (4.88e + 05)	0.002 (0.00)	8.13e + 05 (4.29e + 06)	0.07 (0.03)
IBOSS LASSO	2.17 (1.12)	0.12 (0.02)	3.72 (6.87)	0.18 (0.03)	8.82 (26.93)	0.11 (0.03)	4.91 (11.13)	0.28 (0.05)
SIS	2.03e + 04 (8.24e + 04)	0.002 (0.00)	9.16e + 03 (4.89e + 04)	0.002 (0.002)	10.37 (23.91)	0.003 (0.00)	5.58 (1.26)	0.003 (0.00)
SIS LASSO	1.96 (0.93)	1.18 (0.54)	2.66 (0.69)	0.12 (0.12)	2.81 (0.93)	0.87 (0.39)	3.53 (1.30)	0.13 (0.16)
LASSO	1.83 (0.66)	1.55 (0.62)	1.79 (0.32)	0.26 (0.04)	2.94 (1.80)	1.24 (0.36)	2.70 (0.44)	0.62 (0.07)

Meanwhile, the size of the subdata generated by FAME is the smallest among all. A similar pattern can be observed when 10 features are selected in the subdata.

4.2 | Babycare Process Data in Manufacturing System

The P&G data are offered by Procter & Gamble with 161 continuous observations and 59 continuous features. The response variable is the number of rejects in the production line in a 1-min interval or 5-min interval. By permuting each feature with noise, the number of features increases from 59 to 3009. The subdata

contain 50 observations and 50 features. The estimation result is evaluated using MSE and SNR.

The estimation result is shown in Table 6. For the production line with a 1-min interval, when the subdata are filtered from the original data with 59 features or 3009 features, the proposed FAME method performs the best besides LASSO, then is SIS_LASSO and IBOSS_LASSO. The superiority of FAME is more obvious when the number of features is large. When there are 59 features, the MSE of FAME is 1.87. Increasing the number of features to 3009, the MSE of FAME increases to 2.05 while that of SIS_LASSO increases from 1.96 to 2.66 and that of IBOSS_LASSO increases from 2.17 to 3.72. LASSO always has the smallest MSE because it

TABLE 7 | Results of filtering subdata with 50 observations from five columns of P&G data on prediction variance.

Method	R068_1min						R068_5min					
	5/59			5/3009			5/59			5/3009		
	Prediction variance	True variance	SNR	Prediction variance	True variance	SNR	Prediction variance	True variance	SNR	Prediction variance	True variance	SNR
FAME	0.28	1.80	0.15	0.47	1.80	0.24	0.45	2.67	0.15	0.56	2.67	0.17
IBOSS	2.39e + 07		1.00	4.04e + 08		1.00	1.20e + 05		0.99	7.87e + 05		1.00
IBOSS_LASSO	0.42		0.11	1.92		0.13	6.33		0.22	2.18		0.08
SIS	2.02e + 04		0.80	9.11e + 03		0.50	8.91		0.64	1.35		0.34
SIS_LASSO	0.33		0.10	0.54		0.27	0.29		0.08	0.35		0.11
LASSO	0.21		0.08	0.07		0.04	0.41		0.07	0.04		0.01

uses all original data. However, the proposed method generates subdata with the smallest size and low computational cost.

The superiority of the proposed method compared with benchmark methods, especially for large datasets, can be observed for the production line with 5-min interval. When the number of features is 59, the MSE of FAME is similar to LASSO, but both are larger than the MSE of SIS_LASSO. However, increasing the number of features to 3009, the MSE of FAME is the smallest besides LASSO.

The variance and SNR of the predicted estimator are shown in Table 7. For the production line with a 1-min interval, the variance of the predicted response variable of the proposed FAME method is the smallest besides LASSO. For the production line with a 5-min interval, although the variance of the predicted response variable of FAME is larger than SIS_LASSO, the SNR of FAME is larger than that of SIS_LASSO.

5 | Discussion

In this work, we proposed the FAME method for efficient data reduction for sparse regression. This proposed methodology conducts data reduction in terms of reducing both the number of rows and columns of the original data considering the correlation between predictors and response. Overall, the proposed FAME method is robust in parameter estimation and outperforms benchmark methods in scenarios including continuous response with continuous predictors, continuous response with binary predictors, and binary response with continuous predictors. The superiority is more obvious when the number of data points is much larger than the number of predictors. The advantage of the proposed method comes from selecting the most important features and data points while considering the relationship between the response variable and predictors in the procedure. The subdata generated from the proposed FAME method is always the smallest compared to the benchmark methods.

We would like to remark that the proposed FAME method may not be the best data reduction method but aims to strike a good balance between data reduction efficiency and model estimation accuracy. There are several aspects worth further investigation. Currently, we only discuss the scenarios when the predictors or

response are binary. However, in practice, the non-continuous predictors or response variables usually can have more than two distinct categories. In some cases, both predictors and response can be discrete. This brings the challenge to the proposed FAME algorithm on how to select data points in each level of predictor variables. Besides, for datasets with discrete responses, it is possible that the dataset is highly imbalanced where the majority of classes dominate the whole dataset. It will be of future interest to improve the proposed FAME method to handle imbalanced data. Third, the current FAME method assumes a linear model for data reduction. However, the linear model assumption can be too strong or can be misspecified for some case studies in practice. It will be interesting to extend the proposed method to enable data reduction for nonlinear models or nonparametric models.

Author Contributions

Yanran Wei: conceptualization, methodology, implementation, writing and editing. **William Myers:** conceptualization, editing. **Xinwei Deng:** conceptualization, editing, supervision.

Conflicts of Interest

The authors declare no conflicts of interest.

Data Availability Statement

The data that support the findings of the case study are available from the corresponding author upon reasonable request. Code is available on <https://github.com/Echo-Wei/FAME-Subdata-Selection>. It contains the proposed and benchmark methods as well as examples of implementing those functions in R software.

References

1. P. Drineas, M. W. Mahoney, S. Muthukrishnan, and T. Sarlós, "Faster Least Squares Approximation," *Numerische Mathematik* 117 (2011): 219–249.
2. P. Ma, M. W. Mahoney, and B. Yu, "A Statistical Perspective on Algorithmic Leveraging," *Journal of Machine Learning Research* 16 (2015): 861–911.
3. H. Wang, M. Yang, and J. Stufken, "Information-Based Optimal Subdata Selection for Big Data Linear Regression," *Journal of the American Statistical Association* 114 (2019): 393–405.
4. L. Deldossi and C. Tommasi, "Optimal Design Subsampling From Big Datasets," *Journal of Quality Technology* 54 (2022): 93–101.

5. R. Xie, S. Bai, and P. Ma, "Optimal Sampling Designs for Multi-Dimensional Streaming Time Series With Application to Power Grid Sensor Data," *Annals of Applied Statistics* 17, no. 4 (2023): 3195–3215.

6. C. Meng, R. Xie, A. Mandal, X. Zhang, W. Zhong, and P. Ma, "Low-Con: A Design-Based Subsampling Approach in a Misspecified Linear Model," *Journal of Computational and Graphical Statistics* 30 (2020): 694–708.

7. R. Singh and J. Stufken, "Subdata Selection With a Large Number of Variables," *New England Journal of Statistics in Data Science* 1 (2023): 426–438.

8. Q. Cheng, H. Wang, and M. Yang, "Information-Based Optimal Subdata Selection for Big Data Logistic Regression," *Journal of Statistical Planning and Inference* 209 (2020): 112–122.

9. M. Ai, F. Wang, J. Yu, and H. Zhang, "Optimal Subsampling for Large-Scale Quantile Regression," *Journal of Complexity* 62 (2021): 101512.

10. R. Tibshirani, "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society: Series B: Methodological* 58 (1996): 267–288.

11. J. Fan and R. Li, "Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties," *Journal of the American Statistical Association* 96 (2001): 1348–1360.

12. E. Candes and T. Tao, "The Dantzig Selector: Statistical Estimation When p Is Much Larger Than n ," *Annals of Statistics* 35 (2007): 2313–2351.

13. J. Fan and J. Lv, "Sure Independence Screening for Ultrahigh Dimensional Feature Space," *Journal of the Royal Statistical Society, Series B: Statistical Methodology* 70 (2008): 849–911.

14. H. Wang, "Factor Profiled Sure Independence Screening," *Biometrika* 99 (2012): 15–28.

15. N. Zhao, Q. Xu, M.-L. Tang, B. Jiang, Z. Chen, and H. Wang, "High-Dimensional Variable Screening Under Multicollinearity," *Stat* 9 (2020): e272.

16. B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, "Least Angle Regression," *Annals of Statistics* 32 (2004): 407–499.

17. J. Kiefer, "Optimum Experimental Designs," *Journal of the Royal Statistical Society: Series B: Methodological* 21 (1959): 272–319.

18. R. Killick, P. Fearnhead, and I. A. Eckley, "Optimal Detection of Changepoints With a Linear Computational Cost," *Journal of the American Statistical Association* 107 (2012): 1590–1598.

19. S.-G. Hwang, "Cauchy's Interlace Theorem for Eigenvalues of Hermitian Matrices," *American Mathematical Monthly* 111 (2004): 157–159.

Appendix A

A.1 | Proof of Theorem 1

Proof. For the full data \mathbf{D} with n rows and p columns, denote the subdata obtained after feature screening is $\mathbf{D}_S = (\mathbf{X}_S, \mathbf{y})$ and the one from the FAME algorithm as $\tilde{\mathbf{D}}_S = (\tilde{\mathbf{X}}_S, \tilde{\mathbf{y}})$. $x_{(i)q}$ being the i th order statistic for $\mathbf{x}_q = \{x_{(1)q}, \dots, x_{(n)q}\}$, where $q = 1, \dots, h$ and $i = 1, \dots, n$.

Let $\bar{\mathbf{x}}$ and $\text{var}(\tilde{\mathbf{x}}_q)$ be the mean and variance for the covariate $\tilde{\mathbf{x}}_q$ in $\tilde{\mathbf{D}}_S$. The information matrix of $\tilde{\mathbf{D}}_S$ is

$$\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S = \mathbf{B}^{-1} \begin{bmatrix} k & \mathbf{0}^T \\ \mathbf{0} & (k-1)\mathbf{R} \end{bmatrix} (\mathbf{B}^T)^{-1} \quad (20)$$

where

$$\mathbf{B} = \begin{bmatrix} 1 & & & \\ -\frac{\bar{x}_1}{\sqrt{\text{var}(\tilde{x}_1)}} & \frac{1}{\sqrt{\text{var}(\tilde{x}_1)}} & & \\ & \dots & \ddots & \\ -\frac{\bar{x}_h}{\sqrt{\text{var}(\tilde{x}_h)}} & & & \frac{1}{\sqrt{\text{var}(\tilde{x}_h)}} \end{bmatrix} \quad (21)$$

From Equations (20) and (21),

$$|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S| = k|(k-1)\mathbf{R}| \prod_{q=1}^h \text{var}(\tilde{\mathbf{x}}_q) \geq k(k-1)^h \lambda_{\min}^h(\mathbf{R}) \prod_{q=1}^h \text{var}(\tilde{\mathbf{x}}_q) \quad (22)$$

For each sample variance,

$$(k-1)\text{var}(\tilde{\mathbf{x}}_q) = \sum_{i=1}^k (\tilde{x}_{iq} - \bar{\tilde{x}}_q)^2 \quad (23)$$

$$\geq \left(\sum_{i=1}^l + \sum_{i=n-l+1}^n \right) (x_{(i)q} - \bar{\tilde{x}}_q)^2 \quad (24)$$

$$= \sum_{i=1}^l (x_{(i)q} - \bar{x}_{(u)q})^2 + \sum_{i=n-l+1}^n (x_{(i)q} - \bar{x}_{(v)q})^2 + \frac{l}{2} (\bar{x}_{(u)q} - \bar{x}_{(v)q})^2 \quad (25)$$

$$\geq \frac{l}{2} (\bar{x}_{(u)q} - \bar{x}_{(v)q})^2 \quad (26)$$

$$\geq \frac{l}{2} (x_{(n-l+1)q} - x_{(l)q})^2 \quad (27)$$

where $\bar{\tilde{x}}_q^* = \left(\sum_{i=1}^l + \sum_{i=n-l+1}^n \right) x_{(i)q} / (2l)$, $\bar{x}_{(u)q} = \sum_{i=1}^l x_{(i)q} / l$ and $\bar{x}_{(v)q} = \sum_{i=n-l+1}^n x_{(i)q} / l$. From Equation (27),

$$\text{var}(\tilde{\mathbf{x}}_q) \geq \frac{l(\bar{x}_{(u)q} - \bar{x}_{(v)q})^2}{2(k-1)} \left(\frac{x_{(n-l+1)q} - x_{(l)q}}{\bar{x}_{(u)q} - \bar{x}_{(v)q}} \right)^2 \quad (28)$$

Thus, one can have

$$|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S| \geq k(k-1)^h \lambda_{\min}^h(\mathbf{R}) \prod_{q=1}^h \frac{l(\bar{x}_{(u)q} - \bar{x}_{(v)q})^2}{2(k-1)} \left(\frac{x_{(n-l+1)q} - x_{(l)q}}{\bar{x}_{(u)q} - \bar{x}_{(v)q}} \right)^2 \quad (29)$$

$$= \frac{l^h}{2^h} k \lambda_{\min}^h(\mathbf{R}) \prod_{q=1}^h (\bar{x}_{(u)q} - \bar{x}_{(v)q})^2 \times \prod_{q=1}^h \left(\frac{x_{(n-l+1)q} - x_{(l)q}}{\bar{x}_{(u)q} - \bar{x}_{(v)q}} \right)^2 \quad (30)$$

This shows that

$$\frac{|\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S|}{4^h \prod_{q=1}^h (\bar{x}_{(u)q} - \bar{x}_{(v)q})^2} \geq \frac{\lambda_{\min}^h(\mathbf{R})}{h^h} \times \prod_{q=1}^h \left(\frac{x_{(n-l+1)q} - x_{(l)q}}{\bar{x}_{(u)q} - \bar{x}_{(v)q}} \right)^2 \quad (31)$$

A.2 | Proof of Theorem 2

Proof. Approximations of the lasso estimator of Equation (6) according to Tibshirani [10] is

$$\tilde{\beta}(\lambda) | \tilde{\mathbf{D}}_S \approx \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda) | \tilde{\mathbf{D}}_S) \right]^{-1} \tilde{\mathbf{X}}_S^T \tilde{\mathbf{y}} \quad (32)$$

$$\{\Psi(\tilde{\beta}(\lambda))\} = \text{diag}(\psi_1, \psi_2, \dots, \psi_h) \quad (33)$$

$$\text{where } \psi_j = \begin{cases} \frac{1}{|\tilde{\beta}_j(\lambda)|} & \text{if } \tilde{\beta}_j(\lambda) \neq 0 \\ 0 & \text{otherwise} \end{cases}, \quad \text{for } j = 1, \dots, h \quad (34)$$

The variance of the approximated lasso estimator is

$$\begin{aligned} \text{Var}(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) &\approx \sigma^2 \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1} \tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1} \\ &\quad (35) \end{aligned}$$

Denote β as the OLS estimator of Equation (6) when $\lambda = 0$. Denote $\mathbf{W} = \tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1}$.

$$\begin{aligned} \text{Var}(\beta|\tilde{\mathbf{D}}_S) - \text{Var}(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) &= \sigma^2 (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} - \sigma^2 \mathbf{W}^T (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \mathbf{W} \\ &= \sigma^2 \left[\mathbf{W}^T (\mathbf{W}^T)^{-1} (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \mathbf{W}^{-1} \mathbf{W} - \mathbf{W}^T (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \mathbf{W} \right] \\ &= \sigma^2 \mathbf{W}^T \left[(\mathbf{W}^T)^{-1} (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \mathbf{W}^{-1} - (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \right] \mathbf{W} \\ &= \sigma^2 \mathbf{W}^T \left[(\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right] (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \right. \\ &\quad \left. \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right] (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} - (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \right] \mathbf{W} \\ &= \sigma^2 \mathbf{W}^T \left[2\lambda (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \right. \\ &\quad \left. + \lambda^2 (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \right] \mathbf{W} \\ &= \sigma^2 \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1} \left[2\lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right. \\ &\quad \left. + \lambda^2 \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right] \\ &\quad \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1}. \quad (36) \end{aligned}$$

When $\lambda > 0$, for any $v \neq 0$, $z = \left[\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S + \lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right]^{-1} v \neq 0$.

$$\begin{aligned} v^T (\text{Var}(\beta|\tilde{\mathbf{D}}_S) - \text{Var}(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S)) v &= \sigma^2 z^T \left[2\lambda \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) + \lambda^2 \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) \right] z \\ &= 2\sigma^2 \lambda z^T \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) z + \sigma^2 \lambda^2 z^T \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) (\tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S)^{-1} \Psi(\tilde{\beta}(\lambda)|\tilde{\mathbf{D}}_S) z \\ &> 0. \quad (37) \end{aligned}$$

Equation (36) is positive definite, and it proves that the variance of the lasso estimator decreases as λ increases. Then, from Equations (20) and (21),

$$\text{Var}(\beta|\tilde{\mathbf{D}}_S) = \sigma^2 \tilde{\mathbf{X}}_S^T \tilde{\mathbf{X}}_S^{-1} = \sigma^2 \mathbf{B}^T \begin{bmatrix} \frac{1}{k} & \mathbf{0}^T \\ \mathbf{0} & \frac{1}{k-1} \mathbf{R}^{-1} \end{bmatrix} \mathbf{B} \quad (38)$$

$$\text{Var}(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) \leq V(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) = \frac{\sigma^2}{k-1} \frac{(\mathbf{R}^{-1})_{qq}}{\text{var}(\tilde{x}_q)} \quad (39)$$

where $(\mathbf{R}^{-1})_{qq}$ is the q th diagonal element of \mathbf{R}^{-1} . Denote the spectral decomposition of \mathbf{R} as $\mathbf{R} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$. Since $\Lambda^{-1} \leq \lambda_{\min}^{-1}(\mathbf{R}) \mathbf{I}_h$, $\mathbf{R}^{-1} = \mathbf{V} \mathbf{\Lambda}^{-1} \mathbf{V}^T \leq \mathbf{V} \lambda_{\min}^{-1}(\mathbf{R}) \mathbf{I}_h \mathbf{V}^T = \lambda_{\min}^{-1}(\mathbf{R}) \mathbf{I}_h^T$. Thus, $\mathbf{R}_{qq}^{-1} \leq \lambda_{\min}^{-1}(\mathbf{R})$ for all j . Based on Equations (27) and (38),

$$\begin{aligned} V(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) &\leq \frac{\sigma^2}{k-1} \frac{(\mathbf{R}^{-1})_{qq}}{\text{var}(\tilde{x}_q)} \leq \frac{4h\sigma^2}{k\lambda_{\min}(\mathbf{R})(\bar{x}_{(u)q} - \bar{x}_{(v)q})^2} \\ &\leq \frac{4h\sigma^2}{k\lambda_{\min}(\mathbf{R})(x_{(n-l+1)q} - x_{(l)q})^2} \quad (40) \end{aligned}$$

With the fact that sum of squared differences from the sample mean is smaller than the sum of squared differences from other values,

$$\text{var}(\tilde{x}_q) \leq \frac{1}{k-1} \sum_{i=1}^k \left(\tilde{x}_{iq} - \frac{x_{nq} + x_{1q}}{2} \right)^2 \leq \frac{k}{4(k-1)} (x_{nq} - x_{1q})^2 \quad (41)$$

$$V(\tilde{\beta}_q|\tilde{\mathbf{D}}_S) \geq \frac{4\sigma^2}{k\lambda_{\max}(\mathbf{R})(x_{nq} - x_{1q})^2} \quad (42)$$

In the cases when $\lambda = 0$ and all variables are selected. When $\lambda \rightarrow \infty$, the lower bound of $V(\tilde{\beta}_q|\tilde{\mathbf{D}}_S)$ is close to 0 that none of variables are selected. \square

A.3 | Proof of Theorem 3

We will use the result in Hwang [19] to prove Theorem 3.

Lemma 1 (Cauchy Interlace Theorem). Let \mathbf{A} be a Hermitian matrix of order n and let \mathbf{B} be a principal submatrix of \mathbf{A} of order $n-1$. If $\lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_1$ lists the eigenvalues of \mathbf{A} and $\mu_n \leq \mu_{n-1} \leq \dots \leq \mu_3 \leq \mu_2 \leq \mu_1$ lists the eigenvalues of \mathbf{B} , then $\lambda_n \leq \mu_n \leq \lambda_{n-1} \leq \mu_{n-1} \leq \dots \leq \lambda_2 \leq \mu_2 \leq \lambda_1$.

Proof of Theorem 3. Denote \mathbf{R} the correlation matrix of $\tilde{\mathbf{D}}_S$ selected by the FAME algorithm. \mathbf{R}_{full} as the correlation matrix of the full data \mathbf{D} . Let $\lambda_{\min}(\mathbf{R}) \leq \lambda_{q-1}(\mathbf{R}) \leq \dots \leq \lambda_2(\mathbf{R}) \leq \lambda_{\max}(\mathbf{R})$ lists the eigenvalues of \mathbf{R} , and $\mu_{\min}(\mathbf{R}_{\text{full}}) \leq \mu_{j-1}(\mathbf{R}_{\text{full}}) \leq \dots \leq \mu_2(\mathbf{R}_{\text{full}}) \leq \mu_{\max}(\mathbf{R}_{\text{full}})$ lists the eigenvalues of \mathbf{R}_{full} , where $q = 1, \dots, h$ and $j = 1, \dots, p$. Then

$$\mu_{\min}(\mathbf{R}_{\text{full}}) \leq \lambda_{\min}(\mathbf{R}) \leq \lambda_{\max}(\mathbf{R}) \leq \mu_{\max}(\mathbf{R}_{\text{full}}) \quad (43)$$

Thus, we can have

$$\frac{4\sigma^2}{k\lambda_{\max}(\mathbf{R})(x_{(n)q} - x_{(1)q})^2} \geq \frac{4\sigma^2}{k\mu_{\max}(\mathbf{R}_{\text{full}})(x_{(n)q} - x_{(1)q})^2} \quad (44)$$

$$\frac{4h\sigma^2}{k\lambda_{\min}(\mathbf{R})(\bar{x}_q^u - \bar{x}_q^v)^2} \leq \frac{4h\sigma^2}{k\mu_{\min}(\mathbf{R}_{\text{full}})(\bar{x}_q^u - \bar{x}_q^v)^2} \quad (45) \quad \square$$