

A New Extended BIC and Sequential Lasso Regression Analysis and Their Application in Classification

Jie Chen, Wanzhou Ye

Department of Mathematics, College of Science, Shanghai University, Shanghai, China
Email: listenercjj@shu.edu.cn, wzhy@shu.edu.cn

How to cite this paper: Chen, J. and Ye, W.Z. (2023) A New Extended BIC and Sequential Lasso Regression Analysis and Their Application in Classification. *Advances in Pure Mathematics*, 13, 284-302.
<https://doi.org/10.4236/apm.2023.135020>

Received: April 19, 2023

Accepted: May 28, 2023

Published: May 31, 2023

Copyright © 2023 by author(s) and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

In this paper, firstly, we propose a new method for choosing regularization parameter λ for lasso regression, which differs from traditional method such as multifold cross-validation, our new method gives the maximum value of parameter λ directly. Secondly, by considering another prior form over model space in the Bayes approach, we propose a new extended Bayes information criterion family, and under some mild condition, our new EBIC (NEBIC) is shown to be consistent. Then we apply our new method to choose parameter for sequential lasso regression which selects features by sequentially solving partially penalized least squares problems where the features selected in earlier steps are not penalized in the subsequent steps. Then sequential lasso uses NEBIC as the stopping rule. Finally, we apply our algorithm to identify the nonzero entries of precision matrix for high-dimensional linear discrimination analysis. Simulation results demonstrate that our algorithm has a lower misclassification rate and less computation time than its competing methods under considerations.

Keywords

Regularization Parameter, Sequential Procedure, BIC, Linear Discrimination Analysis, Feature Selection

1. Introduction

Sparse high-dimensional regression (SHR) models arise in many important contemporary scientific fields. A SHR model is:

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

where the number of features p is much larger than the sample size n , and only a relatively small number of the β_j 's are nonzero. Feature selection is crucial in the analysis of SHR models. Desboulets [1] pointed out that there are three main categories of variable selection procedures, they are test-based procedures, penalty-based procedures and screening-based procedures. Regularized regression approaches to the analysis of SHR models have attracted considerable attention of the researchers. A regularized regression approach selects the features and estimates the coefficients simultaneously by minimizing a penalized sum of squares of the form:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \sum_{j=1}^p p_\lambda(|\beta_j|) \quad (2)$$

where λ is a regulating parameter and p_λ is a penalty function such that the number of fitted nonzero coefficients can be regulated by λ ; that is, only a certain number of β_j 's are estimated nonzero when λ is set at a certain value. Various penalty functions have been proposed and studied, including Lasso:

$p_\lambda(|\beta_j|) = \lambda |\beta_j|$, SCAD [2], which smoothly clips a L_1 penalty (for small $|\beta_j|$) and a constant penalty (for large $|\beta_j|$), adaptive Lasso [3]: $p_\lambda(|\beta_j|) = \lambda \omega_j |\beta_j|$, where ω_j are given weights, and MCP [4], which smoothly approaches the L_1 penalty from a constant penalty (for large $|\beta_j|$'s) by an asymptote.

Sequential methods have also received attention in recent decades for feature selection in SHR models. The traditional sequential procedures such as forward stepwise regression (FSR) were criticized for its greedy nature. However, it was discovered recently that the greedy nature is indeed a good one if the goal is to identify relevant features, see [5], especially, in the presence of high spurious correlations due to extremely high dimensionality of the feature space. A sequential procedure of a different nature called least angle regression (LAR) was proposed in [6]. The LAR continuously updates the estimate of the expected responses along a direction having equal angle with the features already selected and selects new features having the largest absolute correlation with the updated current residuals. Recently, Luo and Chen [7] identified the nonzero entries of the precision matrix by a sequential method called JR-SLasso proposed in [8]. Chen and Jiang [9] proposed a two-stage sequential conditional selection (TSCS) approach to the identification and estimation of the nonzeros of the coefficient matrix and the precision matrix. Sequential approach has also been considered for models other than SHR models. Besides their desirable theoretical properties, sequential approaches are computationally more appealing. They are more stable and less affected by the dimensionality of the feature space.

Apart from variable selection, model selection is also an important part for regression analysis. In statistical modeling an investigator often faces the problem of choosing a suitable model from among a collection of viable candidates. Such a determination may be facilitated using a selection criterion, which assigns a score to every model in a candidate set based on some underlying statistical principle. The Bayesian information criterion (BIC) is one of the most widely

known and pervasively used tools in statistical model selection. However, as it was shown by [10] the ordinary Bayesian information criterion is too liberal that is the criteria select far more features than the relevant ones when the model space is large. Chen and Chen [10] proposed a family of extended Bayes information criteria (EBIC) to better meet the needs of variable selection for large model spaces.

Besides regression, classification problem also receives widespread attention in statistical modeling. Classification problems with high-dimensional data rise in many important contemporary scientific fields such as genetics and medicine. The most popular method for classification is Fisher's linear discrimination analysis (LDA). In a K -class classification problem, the LDA assumes that the predictor $x = (x_1, \dots, x_p)^T$ given class $G = k$ follows the multivariate normal distribution $N(u_k, \Sigma)$, $k \in \{1, \dots, K\}$. Let $\pi_k = \Pr(G = k)$ and:

$$d_{kj}(x) = \left[x - \frac{(u_k + u_j)}{2} \right]^T \Omega (u_k - u_j) + \ln(\pi_k / \pi_j), \quad (3)$$

where $\Omega = \Sigma^{-1}$ is the so-called precision matrix. The Bayes rule which is theoretically optimal classifies x into class k if and only if:

$$d_{kj}(x) > 0, \text{ for all } j \neq k.$$

Since the Bayes rule cannot be realized in practice due to the unknown u_k 's and Ω , the Bayes rule is estimated in the LDA by replacing u_k 's and Ω with their estimates. If the dimension of predictor $p (< n)$ is fixed, or diverges under certain conditions, it has been shown that the LDA is asymptotically optimal but that, if $p > n$, it is asymptotically as bad as a random guessing, where n is the sample size, see [11]. The failure of LDA in the case $p > n$ is due to accumulated errors in the estimation of the unknowns, as argued in [12]. Thus, it is necessary to bypass the difficulties in the estimation of unknowns.

In this paper, firstly, we propose a novel method to choose the regularization parameter λ for lasso regression, traditional method such as multifold cross-validation chooses regularization parameter gradually, in details: for lasso regression, the larger regularization parameter λ is, the greater number of regression coefficients will be estimated zero, once λ exceeds the maximum value then there will be no features in the candidate model. Thus, it is important to know the exact maximum value of regularization parameter. However, multifold cross-validation method doesn't tell us the largest value of λ . Differing from multifold cross-validation, our method gives the exact maximum value of regularization parameter such that at least one of the β_j 's will be estimated nonzero. Thus, if you have the data, you can get the largest λ for lasso regression immediately by our method. Then we apply this method to choose parameter for a sequence of partially penalized least squares problems, which means that the features selected in an earlier step are not penalized in the subsequent steps.

Secondly, the re-examination of BIC and EBIC prompts us naturally to consider other forms of prior over the model space in the Bayes approach. We pro-

pose a new extended Bayes information criterion family, and under some mild conditions our new EBIC (NEBIC) is shown to be consistent. Then our NEBIC is used as the stopping rule for sequential lasso algorithm, we dub the proposed procedure as sequential lasso with NEBIC. After that, we apply our algorithm to classification problem, the numerical study demonstrates that our algorithm performs well than its competing methods under consideration.

The remainder of the paper is arranged as follows. In Section 2, we introduce the basic properties of sequential lasso and our new method for choosing the regularization parameter λ . In Section 3, the selection consistency of our NEBIC is established. In Section 4, the structure of our algorithm is given. In Section 5, we apply our algorithm to classification problems, then the main method and simulation results are introduced sequentially.

2. Sequential Lasso Regression for Feature Selection

Let $X = (x_1, x_2, \dots, x_p)$ be the $n \times p$ design matrix, for $j = 1, \dots, p$, $x_j = (x_{1j}, x_{2j}, \dots, x_{nj})^T$ be the observation vector of predictor j on n individuals, $y = (y_1, y_2, \dots, y_n)^T$ be the response vector, and x_j, y all have been standardized, such that $x_j^T 1 = 0$, $y^T 1 = 0$, and $y^T y = n$, $x_j^T x_j = n$. Thus, in model (1) the intercept β_0 can be omitted. Let $\beta = (\beta_1, \beta_2, \dots, \beta_p)^T$, $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^T$, in matrix notation, model (1) is expressed as:

$$y = X\beta + \epsilon \tag{4}$$

Let S denote the set of indices $\{1, 2, \dots, p\}$, let s be any subset of S . Denote by $X(s)$ the matrix consisting of the columns of X with indices in s . Similarly, let $\beta(s)$ denote the vector consisting of the corresponding components of β . Let $\mathcal{R}(s)$ be the linear space spanned by the columns of $X(s)$ and $H(s)$ its corresponding projection matrix:

$$H(s) = X(s)[X^T(s)X(s)]^{-1}X^T(s).$$

At the initial step, sequential lasso minimizes the following penalized sum of squares:

$$l_1 = \|y - X\beta\|_2^2 + \lambda_1 \sum_{j=1}^p |\beta_j| \tag{5}$$

Let $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$, from Theorem 1 we can prove that λ_1 is the largest value of the penalty parameter such that at least one of the β_j^* 's, $j \in \{1, 2, \dots, p\}$, will be estimated nonzero. The features with nonzero estimated coefficients are selected and the set of their indices is denoted by s_{*1} .

For $k \geq 1$, let s_{*k} be the index set of the features selected until step k . At step $k + 1$, sequential lasso minimizes the following penalized sum of squares:

$$l_{k+1} = \|y - X\beta\|^2 + \lambda_{k+1} \sum_{j \in s_{*k}^c} |\beta_j|. \tag{6}$$

From the proposition 1 we can prove that minimization of l_{k+1} is equivalent

to the minimization of

$$\tilde{l}_{k+1} = \|\tilde{y} - \tilde{X}\tilde{\beta}\|_2^2 + \lambda_{k+1} \sum_{j \in S_{*k}^c} |\beta_j|,$$

where $\tilde{y} = [I - H(s_{*k})]y$, $\tilde{X} = [I - H(s_{*k})]X(s_{*k}^c)$, $\tilde{\beta} = \beta(s_{*k}^c)$, $\lambda_{k+1} = 2 \max_{j \in S_{*k}^c} |\tilde{x}_j^T \tilde{y}|$, $\tilde{x}_j = [I - H(s_{*k})]x_j$, from Theorem 1 we can prove that λ_{k+1} is the largest value of the penalty parameter such that at least one of the $\beta_j^2, j \in S_{*k}^c$, will be estimated nonzero.

Next, we give the statements and proofs of proposition 1 and theorem 1.

Proposition 1. For $k \geq 1$, the minimization of l_{k+1} is equivalent to the minimization of \tilde{l}_{k+1} .

Proof: Differentiating l_{k+1} with respect to $\beta(s_{*k})$, we have

$$\frac{\partial l_{k+1}}{\partial \beta(s_{*k})} = -2X^T(s_{*k})y + 2X^T(s_{*k})X(s_{*k})\beta(s_{*k}) + 2X^T(s_{*k})X(s_{*k}^c)\beta(s_{*k}^c).$$

Setting the above derivative to zero, we obtain:

$$\hat{\beta}(s_{*k}) = [X^T(s_{*k})X(s_{*k})]^{-1} X^T(s_{*k})[y - X(s_{*k}^c)\beta(s_{*k}^c)].$$

Substituting $\hat{\beta}(s_{*k})$ into $\|y - X\beta\|^2$ we have:

$$\begin{aligned} l_{k+1} &= \|y - X\beta\|_2^2 + \lambda_{k+1} \sum_{j \in S_{*k}^c} |\beta_j| \\ &= \left\| y - \left(X(s_{*k})\hat{\beta}(s_{*k}) + X(s_{*k}^c)\beta(s_{*k}^c) \right) \right\|_2^2 + \lambda_{k+1} \sum_{j \in S_{*k}^c} |\beta_j| \\ &= \left\| [I - H(s_{*k})][y - X(s_{*k}^c)\beta(s_{*k}^c)] \right\|_2^2 + \lambda_{k+1} \sum_{j \in S_{*k}^c} |\beta_j| \\ &= \|\tilde{y} - \tilde{X}\tilde{\beta}\|_2^2 + \lambda_{k+1} \sum_{j \in S_{*k}^c} |\beta_j| = \tilde{l}_{k+1}. \end{aligned}$$

□

Theorem 1. $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$ is the largest value of the regularization parameter such that at least one of the $\beta_j^2, j \in \{1, 2, \dots, p\}$, will be estimated nonzero.

Proof: this statement is equivalent to that λ_1 is the largest value of the regularization parameter such that the minimization of l_1 will obtain nonzero solutions.

By the KKT condition, let $\frac{\partial l_1}{\partial \beta} = 0$ we have:

$$2X^T(y - X\hat{\beta}) = \lambda_1 \partial \|\hat{\beta}\|_1 \tag{7}$$

In component form:

$$2x_j^T(y - X\hat{\beta}) = \lambda_1 \partial |\hat{\beta}_j|, j \in \{1, 2, \dots, p\}. \tag{8}$$

where $\partial |\hat{\beta}_j|$ is a sub gradient of $|\beta_j|$ at $\hat{\beta}_j$, according to the value of $\hat{\beta}_j$, there are three situations:

$$\begin{aligned} \partial |\hat{\beta}_j| &\in (-1, 1), \text{ when } \hat{\beta}_j = 0 \\ \partial |\hat{\beta}_j| &= 1, \text{ when } \hat{\beta}_j > 0; \\ \partial |\hat{\beta}_j| &= -1, \text{ when } \hat{\beta}_j < 0. \end{aligned}$$

Then (8) is equivalent to:

$$2x_j^T (y - X\hat{\beta}) = \pm\lambda_1, \text{ when } \hat{\beta}_j \neq 0 \tag{9.1}$$

or

$$\left| 2x_j^T (y - X\hat{\beta}) \right| < \lambda_1, \text{ when } \hat{\beta}_j = 0. \tag{9.2}$$

The statement “ λ_1 is the largest value of the penalty parameter such that the minimization of l_1 will obtain nonzero solutions.” is established by showing:

(i) Let $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$, $\min_{\beta} \left\| y - \sum_{j=1}^p \beta_j x_j \right\|_2^2 + \lambda_1 \sum_{j=1}^p |\beta_j|$ will obtain nonzero solutions.

(ii) Let $\lambda > \lambda_1$, $\min_{\beta} \left\| y - \sum_{j=1}^p \beta_j x_j \right\|_2^2 + \lambda \sum_{j=1}^p |\beta_j|$ will only obtain zero solutions.

According to the proof by contradiction, we suppose that there are only zero solutions when $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$.

Substituting $\hat{\beta} = 0$ into (9.2) we have: $\left| 2X^T (y - X * 0) \right| < \lambda_1$, which means:

$$\left| 2x_j^T y \right| < \lambda_1, j \in \{1, 2, \dots, p\}.$$

This contradicts with $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$, thus the supposition is wrong, (i) is proved.

Now let us turn to (ii):

Let $l = \min_{\beta} \left\| y - \sum_{j=1}^p \beta_j x_j \right\|_2^2 + \lambda \sum_{j=1}^p |\beta_j|$, and let $\frac{\partial l}{\partial \beta} = 0$ we have:

$$2X^T (y - X\hat{\beta}) = \lambda \partial \|\hat{\beta}\| \tag{10}$$

Then, the proof of (ii) is equivalent to prove this:

$$\left| 2x_j^T (y - X * 0) \right| < \lambda, j \in \{1, 2, \dots, p\}$$

Substituting $\hat{\beta} = 0$ into the left of (10) we have:

$$\left| 2x_j^T (y - X * 0) \right| = \left| 2x_j^T y \right|, j \in \{1, 2, \dots, p\}$$

By $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$ we have: $\left| 2x_j^T y \right| \leq \lambda_1$, since $\lambda > \lambda_1$, then for all $j \in \{1, 2, \dots, p\}$, we have: $\left| 2x_j^T (y - X * 0) \right| \leq \lambda_1 < \lambda$, thus $\hat{\beta} = 0$ solves equation (10), (ii) is proved. □

3. EBIC and New EBIC

Suppose the dimension of model space S is P , denote S_j is the collection of all models with j covariates, so that the model space S can be partitioned into $S = \bigcup_{j=1}^P S_j$, such that models within each S_j have equal dimension. Let $\tau(S_j)$ be the size of S_j , we know that $\tau(S_j) = C_p^j$. For example, suppose the number of covariates under consideration is $P=1000$, the class of models containing a single covariate is denoted by S_1 , then $S_1 = \{\{1\}, \{2\}, \dots, \{P\}\}$ has size $\tau(S_1) = C_p^1 = 1000$, while the class of models containing two covariate is denoted by S_2 , $S_2 = \{\{1, 2\}, \{1, 3\}, \dots, \{1, P\}, \dots\}$ has size $\tau(S_2) = C_p^2 = (1000 \times 999) / 2$. We can see that the size of S_2 is much bigger than the size of S_1 . Now let us

consider the prior distribution over S as follows.

For $s \in S$, we have $pr(s) = pr(S_j) * pr(s|S_j)$, $j = 1, 2, \dots, P$, since models within each S_j have equal dimension, it is reasonable to assign an equal probability $pr(s|S_j) = \frac{1}{\tau(S_j)}$ for any $s \in S_j$. The ordinary Bayesian information criterion

assigns probabilities $pr(S_j)$ proportional to $\tau(S_j)$, that is $pr(S_j) \propto c * \tau(S_j)$. However, this would cause unreasonable situation by large model space. For example, as we have discussed in the above paragraph, we can see that $\tau(S_2)$ is 999/2 times bigger than $\tau(S_1)$, according to the constant prior by BIC, the probability assigned to S_2 is also 999/2 times that assigned to S_1 . According to the knowledge of combinatorial number, the size of S_j increases as j increases to $j = P/2 = 500$, so that the probability assigned to S_j by the prior increases almost exponentially. In other word, models with larger number of covariates, 50 or 100 say, receive much higher probabilities than models with fewer covariates. This is obviously unreasonable, being strongly against the principle of parsimony.

Instead of assigning probabilities $pr(S_j)$ proportional to $\tau(S_j)$, as in the ordinary BIC, Chen and Chen [10] assign $pr(S_j)$ proportional to $\tau^\varepsilon(S_j)$ for $\varepsilon \in [0, 1]$. This results in the prior probability

$pr(s) = pr(S_j) * pr(s|S_j) = \tau^{\varepsilon-1}(S_j) = \tau^{-\gamma}(S_j)$, $\gamma = 1 - \varepsilon$. This prior distribution gives rise to an extended BIC family (EBIC).

Notice that extended Bayesian information criterion is established by introducing the function x^ε ($x > 1, 0 < \varepsilon < 1$), which aims to select models with fewer covariates. From the perspective of function, we know that x^ε ($x > 1, 0 \leq \varepsilon \leq 1$) is a monotone increasing convex function, and the parameter ε is confined within $[0, 1]$ to ensure upper convex property satisfied. Inspired by this, we consider other upper convex function like $\frac{x}{x+a}$, and the parameter a can be

any positive numbers. In other word, we assign $pr(S_j)$ proportional to $\frac{\tau(S_j)}{\tau(S_j) + a}$, so that

$$pr(s) = pr(S_j) * pr(s|S_j) = \left(\frac{\tau(S_j)}{\tau(S_j) + a} \right) * \frac{1}{\tau(S_j)} = \frac{1}{\tau(S_j) + a}$$

this type of prior distribution on the model space gives rise to a new EBIC family as follows:

$$BIC_a(s) = -2 \ln L_n \{ \hat{\theta}(s) \} + v(s) \ln n + 2 \ln(\tau(S_j) + a), a > 0 \tag{11}$$

where $\hat{\theta}(s)$ is the maximum likelihood estimator of $\theta(s)$ given model s . Now let us investigate the properties of our new EBIC for feature selection in linear regression models. Under some mild conditions our new EBIC is shown to be consistent.

Let y_n be the vector of n observations on the response variable, let X_n be the corresponding design matrix with all the covariates of concerns, and let β be the vector of regression coefficients. Assume that:

$$y_n = X_n \beta + e_n, \tag{12}$$

where $e_n \sim N(0, \sigma^2 I_n)$ and I_n is the identity matrix of size n . Let s_0 be the smallest subset of $\{1, \dots, p_n\}$ such that $u_n = E y_n = X_n(s_0) \beta(s_0)$, where $X_n(s_0)$ and $\beta(s_0)$ are respectively the design matrix and the coefficients corresponding to s_0 . Let $v(s_0)$ be the number of components in s_0 . We call s_0 the true submodel and denote $v(s_0)$ by K_0 , and K is an upper bound for K_0 . Let the projection matrix of $X_n(s)$ be $H_n(s) = X_n(s) [X_n^T(s) X_n(s)]^{-1} X_n^T(s)$. Define $\Delta_n(s) = \|u_n - H_n(s) u_n\|^2$.

The family of our new EBIC under model (12) is defined as:

$$BIC_a(s) = n \ln \left(\frac{y_n^T (I_n - H_n(s)) y_n}{n} \right) + v(s) \ln n + 2 \ln(\tau(S_j) + a), s \in S_j, a > 0. \tag{13}$$

Under the asymptotic identifiability condition proposed by [10], the consistency of our new EBIC is established. The asymptotic identifiability condition is as follows:

Condition 1: asymptotic identifiability. Model (12) with true submodel s_0 is asymptotically identifiable if:

$$\lim_{n \rightarrow \infty} \min \left\{ (\ln n)^{-1} \Delta_n(s) : s \neq s_0, v(s) \leq K_0 \right\} = \infty.$$

And other two lemmas proposed by [13] are also useful to our proof.

Lemma 1: if $\frac{\ln j}{\ln p} \rightarrow \delta$ as $p \rightarrow +\infty$, we have:

$$\ln \left(\frac{p!}{j!(p-j)!} \right) = j \ln p (1 - \delta) (1 + o(1)).$$

Lemma 2: let χ_k^2 denote a χ^2 random variable with degrees of freedom k . If $m \rightarrow +\infty$ and $\frac{K}{m} \rightarrow 0$, then:

$$P(\chi_k^2 \geq m) = \frac{1}{\Gamma(k/2)} (m/2)^{(k/2)-1} e^{-m/2} (1 + o(1)),$$

uniformly for all $k \leq K$.

We now state the consistency result as follows.

Theorem 2. Assume that $p_n = O(n^k)$ for some constant k . If $a > 0$, then under the asymptotic identifiability condition we have:

$$pr \left[\min \{ BIC_a(s) : v(s) = j, s \neq s_0 \} > BIC_a(s_0) \right] \rightarrow 1,$$

for $j = 1, \dots, K$, as $n \rightarrow +\infty$.

Proof:

$$\begin{aligned} & BIC_a(s) - BIC_a(s_0) \\ &= n \ln \left(\frac{y_n^T (I_n - H_n(s)) y_n}{y_n^T (I_n - H_n(s_0)) y_n} \right) + (v(s) - v(s_0)) \ln n, \text{ say,} \\ & \quad + 2 \left[\ln(\tau(S_j) + a) - \ln(\tau(S_{K_0}) + a) \right] \\ & \triangleq T_1 + T_2 \end{aligned}$$

where:

$$\begin{aligned}
 T_1 &= n \ln \left(\frac{y_n^T (I_n - H_n(s)) y_n}{y_n^T (I_n - H_n(s_0)) y_n} \right) \\
 &= n \ln \left(1 + \frac{y_n^T (I_n - H_n(s)) y_n - e_n^T (I_n - H_n(s_0)) e_n}{e_n^T (I_n - H_n(s_0)) e_n} \right), \\
 T_2 &= (v(s) - v(s_0)) \ln n + 2 \left[\ln(\tau(S_j) + a) - \ln(\tau(S_{k_0}) + a) \right].
 \end{aligned}$$

Without loss of generality, we assume that $\sigma^2 = 1$.

Case 1: $s_0 \not\subset s$,

First, we will show $T_1 \geq n \ln \left(1 + \frac{C \ln n}{n} \right)$. We can write:

$$y_n^T (I_n - H_n(s_0)) y_n = e_n^T (I_n - H_n(s_0)) e_n = \sum_{i=1}^{n-v(s_0)} Z_j^2 = n \{1 + o_p(1)\},$$

where Z_j 's are i.i.d. standard normal variables, we have:

$$\begin{aligned}
 &y_n^T (I_n - H_n(s)) y_n - e_n^T (I_n - H_n(s_0)) e_n \\
 &= u_n^T \{I_n - H_n(s)\} u_n + 2u_n^T \{I_n - H_n(s)\} e_n - e_n^T H_n(s) e_n + e_n^T H_n(s_0) e_n \quad (14)
 \end{aligned}$$

By asymptotic identifiability condition, uniformly over s such that $v(s) \leq K$, we have:

$$(\ln n)^{-1} u_n^T \{I_n - H_n(s)\} u_n \rightarrow \infty. \tag{I}$$

Write

$$u_n^T \{I_n - H_n(s)\} e_n = \sqrt{u_n^T \{I_n - H_n(s)\} u_n} Z(s),$$

where:

$$Z(s) = \frac{u_n^T \{I_n - H_n(s)\} e_n}{\sqrt{u_n^T \{I_n - H_n(s)\} u_n}} \sim N(0,1).$$

We hence arrive at

$$\begin{aligned}
 &\max \left[u_n^T \{I_n - H_n(s)\} e_n : s \in S_j \right] \\
 &\leq \sqrt{u_n^T \{I_n - H_n(s)\} u_n} \max \{Z(s) : s \in S_j\} \\
 &\leq \sqrt{u_n^T \{I_n - H_n(s)\} u_n} O_p \left\{ \sqrt{2 \ln p_n} \right\} \\
 &= o_p \left[u_n^T \{I_n - H_n(s)\} u_n \right]
 \end{aligned} \tag{II}$$

where the last inequality follows the Bonferroni inequality, in detail we have:

$$\begin{aligned}
 &P \left(\max \{Z(s) : s \in S_j, j \leq K\} \geq \sqrt{m} \right) \\
 &\leq \sum_{j=1}^K \tau(S_j) P(Z(s) \geq \sqrt{m}) = \sum_{j=1}^K \tau(S_j) P(\chi_1^2 \geq m) \\
 &\leq \sum_{j=1}^K \tau(S_j) P(\chi_j^2 \geq m)
 \end{aligned}$$

where $m = 2K \ln p_n$. The last inequality comes from Lemma 2 that $P(\chi_1^2 \geq m) < P(\chi_j^2 \geq m)$. And according to Lemma 2, we can prove that $\sum_{j=1}^K \tau(S_j) P(\chi_j^2 \geq m) \rightarrow 0$.

Thus, we have

$$\begin{aligned} \max \{e_n^T H_n(s) e_n : v(s) \leq K\} &= m(1 + o_p(1)) \\ &= 2K \ln p_n (1 + o_p(1)) = O_p(\ln p_n) = O_p(\ln n) \end{aligned} \tag{III}$$

We know the term $e_n^T H_n(s_0) e_n$ is a χ^2 -distributed statistic with a fixed degrees of freedom K_0 . Then take (I), (II), (III) into (14) we have:

$$\begin{aligned} &y_n^T (I_n - H_n(s)) y_n - e_n^T (I_n - H_n(s_0)) e_n \\ &= u_n^T \{I_n - H_n(s)\} u_n \left[1 + \frac{2u_n^T \{I_n - H_n(s)\} e_n}{u_n^T \{I_n - H_n(s)\} u_n} - \frac{e_n^T H_n(s) e_n}{u_n^T \{I_n - H_n(s)\} u_n} \right. \\ &\quad \left. + \frac{e_n^T H_n(s_0) e_n}{u_n^T \{I_n - H_n(s)\} u_n} \right] \\ &= u_n^T \{I_n - H_n(s)\} u_n (1 + o_p(1)). \end{aligned}$$

Under the asymptotic identifiability condition, we know that $u_n^T \{I_n - H_n(s)\} u_n$, which goes to infinity faster than $\ln n$, is the dominating term in $y_n^T (I_n - H_n(s)) y_n - e_n^T (I_n - H_n(s_0)) e_n$. Thus,

$$\frac{y_n^T (I_n - H_n(s)) y_n - e_n^T (I_n - H_n(s_0)) e_n}{e_n^T (I_n - H_n(s_0)) e_n} \geq \frac{C \ln n}{n},$$

for any large constant C in probability, thus we have:

$$T_1 = n \ln \left(\frac{y_n^T (I_n - H_n(s)) y_n}{y_n^T (I_n - H_n(s_0)) y_n} \right) \geq n \ln \left(1 + \frac{C \ln n}{n} \right).$$

Next, we will show $BIC_a(s) - BIC_a(s_0) \rightarrow +\infty$, as $n \rightarrow +\infty$.

We know that:

$$\begin{aligned} BIC_a(s) - BIC_a(s_0) &= T_1 + T_2 \\ &\geq n \ln \left(1 + \frac{C \ln n}{n} \right) + (v(s) - v(s_0)) \ln n \\ &\quad + 2 \left[\ln(\tau(S_j) + a) - \ln(\tau(S_{K_0}) + a) \right] \end{aligned} \tag{15}$$

On the one hand,

$$\begin{aligned} &n \ln \left(1 + \frac{C \ln n}{n} \right) + (v(s) - v(s_0)) \ln n \\ &\geq n \ln \left(1 + \frac{C \ln n}{n} \right) - v(s_0) \ln n \\ &= \ln n \frac{n}{\ln n} \left[\ln \left(1 + \frac{C \ln n}{n} \right) - v(s_0) - v(s_0) \frac{\ln n}{n} \right] \end{aligned}$$

Let $\frac{n}{\ln n} = a_n$, we know that $\frac{1}{a_n} \rightarrow 0$ as $n \rightarrow +\infty$, so that $\ln \left(1 + \frac{1}{a_n} \right) \sim \frac{1}{a_n}$.

Thus, the above equation can be expressed as

$$\begin{aligned} &= a_n (\ln n) \left[\ln \left(1 + C \frac{1}{a_n} \right) - K_0 \frac{1}{a_n} \right] \\ &\approx a_n (\ln n) \left[C \frac{1}{a_n} - K_0 \frac{1}{a_n} \right] \\ &= \ln n (C - K_0) > \ln n (C - K) \end{aligned}$$

Which means:

$$n \ln \left(1 + \frac{C \ln n}{n} \right) + (v(s) - v(s_0)) \ln n \geq \ln n (C - K). \tag{IV}$$

On the other hand,

According to Lemma 1 we have $\ln(\tau(S_j) + a) \approx \ln \tau(S_j) \approx j \ln p_n$ as $n \rightarrow +\infty$.

Thus, we obtain:

$$2 \left[\ln(\tau(S_j) + a) - \ln(\tau(S_{K_0}) + a) \right] \approx 2(j - K_0) \ln p_n > -2K \ln p_n. \tag{V}$$

Bring (IV), (V) into (15) we can see that choosing $C > K(2k + 1)$, we obtain:

$$BIC_a(s) - BIC_a(s_0) > \ln n (C - K(2k + 1)) \rightarrow +\infty.$$

Case 2: $s_0 \subset s$,

First, we will show:

$$T_1 = n \ln \left(\frac{y_n^T (I_n - H_n(s)) y_n}{y_n^T (I_n - H_n(s_0)) y_n} \right) \geq -2j \ln p_n \{1 + o_p(1)\}.$$

In this case we have $\{I_n - H_n(s)\} X_n(s_0) = 0$, thus $y_n^T (I_n - H_n(s)) y_n = e_n^T (I_n - H_n(s)) e_n$.

And

$$\begin{aligned} &e_n^T (I_n - H_n(s_0)) e_n - y_n^T (I_n - H_n(s)) y_n \\ &= e_n^T (H_n(s) - H_n(s_0)) e_n = \sum_{i=1}^j Z_i^2(s). \end{aligned}$$

where $j = v(s) - v(s_0)$, $Z_i(s)$ are some independent standard normal random variables depending on s . Let $\hat{e}_n = \{I_n - H_n(s_0)\} e_n$, we obtain that:

$$\begin{aligned} &n \ln \left(\frac{y_n^T (I_n - H_n(s_0)) y_n}{y_n^T (I_n - H_n(s)) y_n} \right) \\ &= n \ln \left(1 + \frac{e_n^T (I_n - H_n(s_0)) e_n - e_n^T (I_n - H_n(s)) e_n}{e_n^T (I_n - H_n(s)) e_n} \right) \\ &= n \ln \left\{ 1 + \frac{\sum_{i=1}^j Z_i^2(s)}{\hat{e}_n^T \hat{e}_n - \sum_{i=1}^j Z_i^2(s)} \right\} \leq n \frac{\sum_{i=1}^j Z_i^2(s)}{\hat{e}_n^T \hat{e}_n - \sum_{i=1}^j Z_i^2(s)}. \end{aligned}$$

As $n \rightarrow +\infty$, $n^{-1} \hat{e}_n^T \hat{e}_n \rightarrow \sigma^2 = 1$. And in case 1 we have shown that:

$$\max \left\{ \sum_{i=1}^j Z_i^2(s), s \in S_{v(s)} \right\} = 2j \ln p_n \{1 + o_p(1)\}.$$

Thus,

$$\begin{aligned} & \max \left\{ \frac{n \sum_{i=1}^j Z_i^2(s)}{\hat{e}_n^T \hat{e}_n - \sum_{i=1}^j Z_i^2(s)} : s \in S_{v(s)} \right\} \\ & \leq \frac{2nj \ln p_n \{1 + o_p(1)\}}{(n - 2j \ln p_n) \{1 + o_p(1)\}} = 2j \ln p_n \{1 + o_p(1)\}. \end{aligned}$$

So, we can see that:

$$T_1 = -n \ln \left(\frac{y_n^T (I_n - H_n(s_0)) y_n}{y_n^T (I_n - H_n(s)) y_n} \right) \geq -2j \ln p_n \{1 + o_p(1)\}.$$

Consequently, uniformly in s such that $v(s) = j + v(s_0)$, we have:

$$\begin{aligned} BIC_a(s) - BIC_a(s_0) &= T_1 + T_2 \\ &\geq -2j \ln p_n \{1 + o_p(1)\} + (v(s) - v(s_0)) \ln n \\ &\quad + 2 \left[\ln(\tau(S_j) + a) - \ln(\tau(S_{K_0}) + a) \right] \\ &\approx -2j \ln p_n \{1 + o_p(1)\} + (v(s) - v(s_0)) \ln n + 2(v(s) - v(s_0)) \ln p_n \\ &= -2j \ln p_n \{1 + o_p(1)\} + j \ln n + 2j \ln p_n = j \ln n. \end{aligned}$$

Thus, we have:

$BIC_a(s) - BIC_a(s_0) \geq j \ln n \rightarrow +\infty$ as $n \rightarrow +\infty$. The conclusion hence follows. □

4. The Structure of Our Algorithm

Initial step:

Standardize $y, x_j, j = 1, \dots, p$, such that $y^T 1 = 0, x_j^T 1 = 0, y^T y = n, x_j^T x_j = n$, let $\lambda_1 = 2 \max_{j \in \{1, 2, \dots, p\}} |x_j^T y|$, consider the following penalized sum of squares:

$$\min_{\beta} \left\{ \|y - X\beta\|_2^2 + \lambda_1 \|\beta\|_1 \right\},$$

the features with nonzero estimated coefficients are selected and the set of their indices is denoted by S_{TEMP} , let $s_{*1} = S_{TEMP}$. Compute $I - H(s_{*1})$ and $NEBIC(s_{*1})$.

General step:

For $k \geq 1$, compute $\tilde{x}_j \tilde{y}$ for $j \in S_{*k}^c$, where $\tilde{y} = [I - H(s_{*k})]y, \tilde{x}_j = [I - H(s_{*k})]x_j$. Let $\lambda_{k+1} = 2 \max_{j \in S_{*k}^c} |\tilde{x}_j \tilde{y}|$, consider the following penalized sum of squares:

$$\min_{\beta} \left\{ \|\tilde{y} - \tilde{X}\tilde{\beta}\|_2^2 + \lambda_{k+1} \|\tilde{\beta}\|_1 \right\},$$

the features with nonzero estimated coefficients are selected and the set of their indices is denoted by S_{TEMP} , let $s_{*k+1} = S_{*k} \cup S_{TEMP}$. Compute $NEBIC(s_{*k+1})$, if $NEBIC(s_{*k+1}) > NEBIC(s_{*k})$, stop; otherwise, compute $I - H(s_{*k+1})$ and continue.

Suppose that the above procedure stops at $k^* + 1$ step, then our algorithm outputs the index set of selected features which is denoted by s_{*k^*} . Then the parameters in the selected model are estimated by their least-square estimates:

$$\hat{\beta}(s_{*k^*}) = [X^T(s_{*k^*})X(s_{*k^*})]^{-1} X^T(s_{*k^*})y, \hat{\beta}(s_{*k^*}^c) = 0.$$

The NEBIC for s_{*k} , $k = 1, 2, \dots, k^*$ in the above algorithm is given by:

$$\begin{aligned} \text{NEBIC}(s_{*k}) = n \ln \left(\frac{y^T (I_n - H_n(s_{*k}))y}{n} \right) + \nu(s_{*k}) \ln n \\ + 2 \ln \left(\tau(S_{\nu(s_{*k})}) + a \right) \text{ where } a > 0. \end{aligned}$$

5. Application to High-Dimensional Classification

5.1. Method

As we have discussed before, it is necessary to give better methods to estimate the unknown u_k 's and Ω for the small- n -large- p problem. In this paper, we apply our algorithm to identify the non-zero entries of precision matrix, and the estimate of Ω is then obtained by the constrained maximum likelihood estimation. And we adopt the method proposed by [7] to estimate class means. The estimated class means and precision matrix are finally used to construct the discrimination rule.

5.2. Procedure

1) Constrained estimation for the class means

The predictor components are first ordered according to the F-statistic for testing the significance of class effects, instead of being pairwise fused, the class means for each component are then clustered by a divisive hierarchical clustering procedure. The class means are estimated under the structure revealed by the clustering. For details, we refer the reader to [7].

2) Constrained estimation for the precision matrix

The identification of non-zero entries in a concentration matrix has attracted considerable attention of the researchers, concentration matrix is the inverse of the covariance matrix of a random vector. A concentration matrix is closely related to an undirected graphical model. An undirected graphical model is specified by a vertex set V and an edge set E , and is denoted by $G = (V, E)$. The vertex set V represents a collection of random variables $\{Y_1, \dots, Y_p\}$. The edge set E describes the inter-relationship among the random variables: there is an edge connecting vertices Y_i and Y_j if they are dependent conditioning on all the remaining variables. Suppose that $Y = (Y_1, \dots, Y_p)$ follows a multivariate normal distribution with concentration matrix $\Omega = (\Omega_{ij})$. Then, there is an edge between Y_i and Y_j if and only if $\Omega_{ij} = \Omega_{ji} \neq 0$. Thus, the detection of edges of G is equivalent to the identification of non-zero entries of Ω .

In this paper, we adopt the method, which was proposed by [14], is based on

the relationship between entries of Ω and the coefficients of p regression models where each component of Y is regressed on the remaining $p - 1$ components. A non-zero entry of Ω corresponds to a non-zero regression coefficient in the regression models. In other word, the detection and estimation of non-zero entries of Ω are then boiled down to the selection and estimation of non-zero coefficients in p regression models. According to this, we first apply our algorithm to identify non-zero entries of Ω , and the estimate of Ω is then obtained by the constrained maximum likelihood estimation. The details as follow.

For an undirected graph $G = (V, E)$, let V be modeled as the set of the components of a random vector $Y = (Y_1, \dots, Y_p)$. We assume that Y follows a multivariate normal distribution $N(u, \Sigma)$. Without loss of generality, assume that $u = 0$. Let Y_{i-} be the vector obtained from Y by eliminating component Y_i . Denote by Σ_{i-i-} the variance-covariance matrix of Y_{i-} , by σ_i^2 the variance of Y_i , and by Σ_{ii-} the covariance vector between Y_i and Y_{i-} . We know that the conditional distribution of Y_i given Y_{i-} is still normal, with the following conditional mean and conditional variance:

$$E(Y_i | Y_{i-}) = \Sigma_{ii-} \Sigma_{i-i-}^{-1} Y_{i-}, \text{Var}(Y_i | Y_{i-}) = \sigma_i^2 - \Sigma_{ii-} \Sigma_{i-i-}^{-1} \Sigma_{i-i-}$$

Let β_{ij} be the j th component of $\Sigma_{ii-} \Sigma_{i-i-}^{-1}$. We can then express the conditional distributions in the form of the following regression models:

$$Y_i = \sum_{j \neq i} \beta_{ij} Y_j + \epsilon_i, \epsilon_i \sim N(0, D_i), i = 1, \dots, p, \tag{16}$$

where $D_i = \sigma_i^2 - \Sigma_{ii-} (\Sigma_{i-i-})^{-1} \Sigma_{i-i-}$. Without loss of generality, suppose that Y_i is the first component of Y . Let the covariance matrix Σ be partitioned as:

$$\Sigma = \begin{pmatrix} \sigma_i^2 & \Sigma_{ii-} \\ \Sigma_{i-i} & \Sigma_{i-i-} \end{pmatrix}, \text{let } \Omega = (\omega_{jl})_{j,l \in \{1, \dots, p\}}, \text{ we can obtain:}$$

$$\Omega = \Sigma^{-1} = \begin{pmatrix} D_i^{-1} & -D_i^{-1} \Sigma_{ii-} \Sigma_{i-i-}^{-1} \\ -\Sigma_{i-i-}^{-1} \Sigma_{i-i} D_i^{-1} & \Sigma_{i-i-}^{-1} + \Sigma_{i-i-}^{-1} \Sigma_{i-i} D_i^{-1} \Sigma_{ii-} \Sigma_{i-i-}^{-1} \end{pmatrix}.$$

By comparing the left upper block of Ω with $\Sigma_{ii-} \Sigma_{i-i-}^{-1}$, we see that:

$$\beta_{ij} = -\frac{\omega_{ij}}{D_i^{-1}} = -\frac{\omega_{ij}}{\omega_{ii}},$$

noting that $D_i^{-1} = \Omega_{ii}$. These connections establish the equivalence:

$$\omega_{ij} = 0 \Leftrightarrow \beta_{ij} = 0.$$

As we can see that the identification of non-zero entries of Ω reduces to the identification of the non-zero β_{ij} in the above regression models. Thus, firstly, we apply our algorithm to model (16) to identify all the non-zero entries of Ω . Let $\mathcal{E} = \{(j, l) : \omega_{jl} \neq 0\}$, according to our algorithm we have $\hat{\mathcal{E}}$.

Next, we summary our proposed approaches. First, we introduce some notation. Denote by y_i the standardized n -vector of observed Y_i , let X_i be the $n \times (p - 1)$ matrix consisting of all y_j with $j \neq i$. Here is our algorithm based on above data sets:

Algorithm for identifying \mathcal{E}

Initial step:

Let $\lambda_i = 2 \max |X_i^T y_i|$, let β_i be $(p-1)$ -vector, for $i=1,2,\dots,p$. Consider minimizing following p penalized sum of squares separately:

$$\min_{\beta_i} \frac{1}{2} \|y_i - X_i \beta_i\|_2^2 + \lambda_i \|\beta_i\|_1, i=1,2,\dots,p,$$

Let s_{1*1}, \dots, s_{p*1} denote the indices of features with nonzero estimated coefficients for the above p regression models respectively. Let $s_{*1} = \{s_{1*1}, \dots, s_{p*1}\}$. Then compute $\text{NEBIC}(s_{*1})$ and $I - H(s_{*1})$ for $i=1,2,\dots,p$.

General step $k + 1$ ($k \geq 1$):

Let $\lambda_i = 2 \max |X_i^T (I - H(s_{i*k})) y_i|$ for $i=1,2,\dots,p$. Consider minimizing following p penalized sum of squares separately:

$$\min_{\tilde{\beta}_i} \frac{1}{2} \|\tilde{y}_i - \tilde{X}_i \tilde{\beta}_i\|_2^2 + \lambda_i \|\tilde{\beta}_i\|_1, i=1,2,\dots,p,$$

where $\tilde{y}_i = (I - H(s_{i*k})) y_i$, $\tilde{X}_i = (I - H(s_{i*k})) X_i$, $\tilde{\beta}_i = \beta_i(s_{i*k}^c)$. Let s_{iTEMP} be the indices of features with nonzero estimated coefficients for the above i th regression model, and $s_{i*k+1} = s_{iTEMP} \cup s_{i*k}$ for $i=1,2,\dots,p$. Then let $s_{*k+1} = \{s_{1*k+1}, \dots, s_{p*k+1}\}$. Compute $\text{NEBIC}(s_{*k+1})$. If $\text{NEBIC}(s_{*k+1}) > \text{NEBIC}(s_{*k})$, stop; otherwise, compute $I - H(s_{i*k+1})$ for $i=1,2,\dots,p$ and continue.

Suppose the above algorithm stops at step $k^* + 1$, then our algorithm output the index set $s_{*k^*} = \{s_{1*k^*}, \dots, s_{p*k^*}\}$.

The NEBIC for set s_{*k} is given by the following formula:

$$\text{NEBIC}(s_{*k}) = \sum_{i=1}^p \left\{ n \ln \left(\left\| [I - H(s_{i*k})] y_i \right\|_2^2 \right) + |s_{i*k}| \ln n + 2 \ln \left(\tau(S_{|s_{i*k}|}) + a \right) \right\}.$$

Estimation of precision matrix:

From the above algorithm we identify all the non-zero entries of $\hat{\Omega}$, let $\hat{\mathcal{E}}$ be the set of all the indices. Then $\hat{\Omega}$ is obtained by the constrained maximum likelihood estimation as follows:

$$\hat{\Omega} = \min_{\omega_{jl}=0:(j,l) \in \hat{\mathcal{E}}} \left\{ \text{trace}(\hat{S}\Omega) - \ln \det(\Omega) + \lambda \|\Omega\|_1 \right\},$$

where \hat{S} is the empirical covariance matrix.

3) Discrimination rule

The discrimination rule is constructed by replacing the unknowns in $d_{kj}(x)$ by their estimates obtained above. Here $d_{kj}(x)$ is:

$$d_{kj}(x) = \left[x - (u_k + u_j) / 2 \right]^T \Omega (u_k - u_j) + \ln(\pi_k / \pi_j).$$

Then we have $\hat{d}_{kj}(x)$:

$$\hat{d}_{kj}(x) = \left[x - (\hat{u}^k + \hat{u}^j) / 2 \right]^T \hat{\Omega} (\hat{u}^k - \hat{u}^j) + \ln(n_k / n_j).$$

Classify x into class k if and only if:

$$\hat{d}_{kj}(x) > 0, \text{ for all } j \neq k.$$

5.3. Simulation Studies

We compare our method with other methods available in the literature through numerical studies in this section. The methods considered for the comparison are: 1) linear discrimination with detected sparsity (LDwDS) in [7]. 2) The adaptive hierarchically penalized nearest shrunken centroid (ahp-NSC) and the adaptive L_∞ -norm penalized NSC (alp-NSC) proposed in [15]. 3) Pairwise SIS (PSIS) in [16].

The performance of the methods is evaluated by the misclassification rate (MCR) and computation time. We have three simulation settings, repeat each setting 100 times. At each replicate under a simulation setting, we simulate a training data set and a testing data set. The training data set is used to construct discrimination rules with the methods, and the testing set is used to evaluate the MCR and record computation time.

We design the simulation setting inspired by [7], consider $K = 2$, $p = 400$. The following single scheme for the class means is taken throughout:

$$u_1 = 0, u_2 = \left(\underbrace{0.5, \dots, 0.5}_{100}, \underbrace{0, \dots, 0}_{p-100} \right).$$

The covariance matrix is generated through the precision matrix. First, the nonzero positions of the precision matrix are decided in the following schemes:

ES1: the non-zero entries of Ω is randomly determined with:

$$\Pr(\omega_{ij} = 0) = 1 - \Pr(\omega_{ij} \neq 0) = 0.99.$$

ES2: Ω is a diagonal block matrix with 10 blocks of size 10.

ES3: Ω is a diagonal block matrix with 10 blocks of size 10, each block is a diagonal band matrix with $\omega_{ij} \neq 0$ if $|i - j| \leq 2$.

For the nonzero values of the precision matrix, we first generate $\tilde{\Omega} = \tilde{\omega}_{ij}$ as follows: $\tilde{\omega}_{ii} = 1$, $\tilde{\omega}_{ij} = \tilde{\omega}_{ji}$ are generated as i.i.d. observations from the uniform distribution on $[-0.3, 0.7]$. Then we take $\Omega = \tilde{\Omega} + (0.1 - \lambda_{\min}(\tilde{\Omega}))I$. Eventually, take the common covariance matrix as the correlation matrix corresponding to Ω^{-1} . The training sample size $n_1 = 200$ and the testing sample size $n_2 = 1000$. The covariance matrix is generated as a diagonal block matrix with four 100×100 identical blocks. The 100×100 block is generated by the same generating schemes above. The simulation results under these three settings are reported in **Tables 1-3**.

From all three tables we can see that as for the MCR, the performance of our method is much better than the two NSC methods and PSIS. Although our method has the same MCR with LDwDS from **Table 1** and **Table 2**, it takes shorter computation time than LDwDS. From **Table 3** we find that our method performs better than LDwDS both in misclassification and computation time. From [7] we notice that the misclassification rates of LDwDS are universally lower than all the other methods under consideration. Nevertheless, it is not without drawbacks. LDwDS takes longer computation time than other methods.

Table 1. Averaged misclassification rate (MCR) and computation time (in seconds) under simulation ES1 (numbers in the parentheses are standard deviations).

Method	MCR	Time
LDwDS	0.005 (0.005)	14.339 (1.502)
ahp-NSC	0.045 (0.026)	0.213 (0.049)
alp-NSC	0.193 (0.050)	0.202 (0.026)
PSIS	0.500 (0.006)	0.115 (0.009)
our method	0.005 (0.005)	11.008 (0.962)

Table 2. Averaged misclassification rate (MCR) and computation time (in seconds) under simulation ES2 (numbers in the parentheses are standard deviations).

Method	MCR	Time
LDwDS	0.001 (0.003)	28.532 (3.524)
ahp-NSC	0.147 (0.043)	0.205 (0.027)
alp-NSC	0.291 (0.043)	0.199 (0.024)
PSIS	0.500 (0.006)	0.114 (0.007)
our method	0.001 (0.002)	10.684 (1.199)

Table 3. Averaged misclassification rate (MCR) and computation time (in seconds) under simulation ES3 (numbers in the parentheses are standard deviations).

Method	MCR	Time
LDwDS	0.017 (0.025)	20.852 (2.384)
ahp-NSC	0.167 (0.049)	0.194 (0.007)
alp-NSC	0.310 (0.041)	0.190 (0.009)
PSIS	0.500 (0.006)	0.110 (0.004)
our method	0.016 (0.023)	8.725 (0.709)

Therefore, from the perspectives of computation time and misclassification, our algorithm performances better than its competing algorithms.

6. Summary and Discussion

In this paper, on the one hand, based on the characteristic of lasso regression we propose a novel method to choose the regularization parameter λ , which gives the exact maximum value of λ for lasso regression. On the other hand, since the ordinary Bayes information criterion is too liberal for model selection when the model space is large, inspired by [10] we propose a new EBIC (NEBIC). Compared with EBIC the parameter in our new criterion doesn't need to be restricted in a small range, what's more, under some mild conditions our new EBIC is also shown to be consistent. Then, based on these two findings, we have a new algorithm for feature selection. In detail, we apply our new method to choose regula-

rization parameter for sequential lasso, and our NEBIC is used as the stopping rule. After that, we apply our algorithm to classification problem, the numerical study demonstrates that our algorithm performs better than its competing methods under consideration.

Further research may consider these two questions:

Firstly, our method chooses the largest regularization parameter λ such that at least one of the β_j 's, $j \in \{1, 2, \dots, p\}$, will be estimated nonzero. A nature question is that instead of solving the optimization problems can we obtain the indices of those features with nonzero estimated coefficients by easier methods? Luo and Chen [17] pointed out that these indices are related to the choice of regularization parameter λ under some conditions. However, these conditions are too strict. Further research may consider if there exist milder conditions to connect our new method for choosing regularization parameter with the indices of non-zero estimated coefficients.

Secondly, we choose the upper convex function $\frac{x}{x+a}$, $a > 0$ as the prior function over the model space to extend EBIC, and our new EBIC is shown to be consistent. Subsequent research may consider if there exists a special function class for model space priors such that within this class the Bayes information criteria is consistent.

Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

References

- [1] Desboulets, L.D.D. (2018) A Review on Variable Selection in Regression Analysis. *Econometrics*, **6**, Article 45. <https://doi.org/10.3390/econometrics6040045>
- [2] Fan, J. and Li, R. (2001) Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties. *Journal of the American Statistical Association*, **96**, 1348-1360. <https://doi.org/10.1198/016214501753382273>
- [3] Zou, H. (2012) The Adaptive Lasso and Its Oracle Properties. *Journal of the American Statistical Association*, **101**, 1418-1429. <https://doi.org/10.1198/016214506000000735>
- [4] Zhang, C. (2010) Nearly Unbiased Variable Selection under Minimax Concave Penalty. *The Annals of Statistics*, **38**, 894-942. <https://doi.org/10.1214/09-AOS729>
- [5] Tropp, J.A. (2004) Greed Is Good: Algorithmic Results for Sparse Approximation. *IEEE Transactions on Information Theory*, **50**, 2231-2242. <https://doi.org/10.1109/TIT.2004.834793>
- [6] Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004) Least Angle Regression. *The Annals of Statistics*, **32**, 407-499. <https://doi.org/10.1214/009053604000000067>
- [7] Luo, S. and Chen, Z. (2020) A Procedure of Linear Discrimination Analysis with Detected Sparsity Structure for High-Dimensional Multi-Class Classification. *Journal of Multivariate Analysis*, **179**, Article ID: 104641. <https://doi.org/10.1016/j.jmva.2020.104641>
- [8] Luo, S. and Chen, Z. (2014) Edge Detection in Sparse Gaussian Graphical Models.

Computational Statistics & Data Analysis, **70**, 138-152.

<https://doi.org/10.1016/j.csda.2013.09.002>

- [9] Chen, Z. and Jiang, Y. (2020) A Two-Stage Sequential Conditional Selection Approach to Sparse High-Dimensional Multivariate Regression Models. *Annals of the Institute of Statistical Mathematics*, **72**, 65-90.
<https://doi.org/10.1007/s10463-018-0686-5>
- [10] Chen, J. and Chen, Z. (2008) Extended Bayesian Information Criteria for Model Selection with Large Model Spaces. *Biometrika*, **95**, 759-771.
<https://doi.org/10.1093/biomet/asn034>
- [11] Bickel, P.J. and Levina, E. (2004) Some Theory for Fisher's Linear Discriminant Function, 'Naive Bayes', and Some Alternatives When There Are Many More Variables than Observations. *Bernoulli*, **10**, 989-1010.
<https://doi.org/10.3150/bj/1106314847>
- [12] Fan, J. and Fan, Y. (2008) High Dimensional Classification Using Features Annealed Independence Rules. *Annals of Statistics*, **36**, 2605-2637.
<https://doi.org/10.1214/07-AOS504>
- [13] Luo, S. and Chen, Z. (2013) Extended BIC for Linear Regression Models with Diverging Number of Relevant Features and High or Ultra-High Feature Spaces. *Journal of Statistical Planning and Inference*, **143**, 494-504.
<https://doi.org/10.1016/j.jspi.2012.08.015>
- [14] Meinshausen, N. and Bühlmann, P. (2006) High-Dimensional Graphs and Variable Selection with the Lasso. *The Annals of Statistics*, **34**, 1436-1462.
<https://doi.org/10.1214/009053606000000281>
- [15] Wang, S. and Zhu, J. (2007) Improved Centroids Estimation for the Nearest Shrunken Centroid Classifier. *Bioinformatics*, **23**, 972-979.
<https://doi.org/10.1093/bioinformatics/btm046>
- [16] Pan, R., Wang, H. and Li, R. (2016) Ultrahigh-Dimensional Multiclass Linear Discriminant Analysis by Pairwise Sure Independence Screening. *Journal of the American Statistical Association*, **111**, 169-179.
<https://doi.org/10.1080/01621459.2014.998760>
- [17] Luo, S. and Chen, Z. (2014) Sequential Lasso Cum EBIC for Feature Selection with Ultra-High Dimensional Feature Space. *Journal of the American Statistical Association*, **109**, 1229-1240. <https://doi.org/10.1080/01621459.2013.877275>