

Minimization of Akaike's Information Criterion in Linear Regression Analysis via Mixed Integer Nonlinear Program

Keiji Kimura^{*1} and Hayato Waki^{†2}

¹Faculty of Mathematics, Kyushu University

²Institute of Mathematics for Industry, Kyushu University

First version : June 17, 2016, Revised : June 24, 2016

Akaike's information criterion (AIC) is a measure of the quality of a statistical model for a given set of data. We can determine the best statistical model for a particular data set by the minimization of the AIC. Since we need to evaluate exponentially many candidates of the model by the minimization of the AIC, the minimization is unreasonable. Instead, stepwise methods, which are local search algorithms, are commonly used to find a better statistical model though it may not be the best.

We propose a branch and bound search algorithm for a mixed integer nonlinear programming formulation of the AIC minimization by Miyashiro and Takano (2015). More concretely, we propose methods to find lower and upper bounds, and branching rules for this minimization. We then combine them with SCIP, which is a mathematical optimization software and a branch-and-bound framework. We show that the proposed method can provide the best statistical model based on AIC for small-sized or medium-sized benchmark data sets in UCI Machine Learning Repository. Furthermore, we show that this method finds good quality solutions for large-sized benchmark data sets.

Keywords : Mixed integer nonlinear program, branch-and-bound, SCIP and Akaike's information criterion

1 Introduction

Selecting the best statistical model from a number of candidate statistical models for a given set of data is one of the most important problems solved in statistical applications, *e.g.* regression analysis. This is called *variable selection*. The purposes of variable selection are to provide the

^{*}744 Motooka, Nishi-ku, Fukuoka 819-0395, Japan. k-kimura@math.kyushu-u.ac.jp

[†]744 Motooka, Nishi-ku, Fukuoka 819-0395, Japan. waki@imi.kyushu-u.ac.jp

simplest statistical model for a given data set and to improve the prediction performance while keeping the goodness-of-fit for a given data set. See [8] for more details on variable selection.

In variable selection based on *an information criterion*, all the candidates are evaluated by the information criterion and select a statistical model by using those evaluations. Akaike's information criterion (AIC) is one of the information criteria and proposed in [3]. An AIC value is computed for each candidate, and the model whose AIC value is the smallest is selected as the best statistical model. Since we often need to handle too many candidates of statistical models in practical applications, the global minimization based on AIC is not practical. Instead of the global minimization, stepwise methods, which are local search algorithms, are commonly used to find a statistical model which has as small AIC as possible, but it may not be the smallest.

The contribution of our study is to propose a branch and bound search algorithm for a mixed integer nonlinear programming (MINLP) formulation of the minimization of AIC in linear regression by Miyashiro and Takano [12]. Miyashiro and Takano [12] propose a mixed integer second-order cone programming (MISOCP) formulation from the MINLP formulation and solve the resulting problems by CPLEX [9], while we propose procedures to find lower and upper bounds of the MINLP problems and define branching rules for efficient computation. In addition, we provide an implementation to solve it efficiently via SCIP. SCIP is a mathematical optimization software and a branch-and-bound framework. SCIP has high flexibility of user plugin and control on various parameters in the branch-and-bound framework for efficient computation. We also propose an efficient computation for a set of data which has linear dependency. By applying our proposed method to benchmark data sets in [16], we can obtain the best statistical models for some of them. Our implementation is available at [18].

We introduce some related work. Miyashiro and Takano [12] propose a MISOCP formulation for variable selection based on some information criteria in linear regression. Bertsimas and Shioda [6] and Bertsimas, King and Mazumder [5] provide a mixed integer quadratic programming (MIQP) formulation for linear regression with a cardinality constraint. Their formulation is available to our problems by fixing the number of explanatory variables. We compare our proposed method with MIQP and MISOCP formulations, and observe that our proposed method outperforms MIQP and MISOCP formulations.

The organization of this manuscript is as follows: We give a brief introduction of linear regression based on AIC in Section 2. We introduce the MINLP formulation of the AIC minimization and ways to find lower and upper bounds used in the branch-and-bound framework in Section 3. Section 4 introduces techniques for more efficient computation, *e.g.* branching rules and treatment on data which has linear dependency. We present numerical results in Section 5. In particular, we show the numerical comparison with MISOCP and MIQP formulations. In addition, we present numerical performances of branching rules proposed in subsection 4.4. We discuss future work of our proposed method in Section 6. This manuscript is a full paper version of [10].

2 Preliminary on Akaike's information criterion in linear regression

We explain how to select the best statistical model via AIC in linear regression analysis. Linear regression is a fundamental statistical tool which determines coefficients $\beta_0, \dots, \beta_p \in \mathbb{R}$ for the

following equation from a given set of data:

$$y = \beta_0 + \sum_{j=1}^p \beta_j x_j. \quad (1)$$

Here x_1, \dots, x_p and y are called *the explanatory variables* and *the response variable* respectively. In fact, we adopt coefficients β_0, \dots, β_p which minimize $\sum_{i=1}^n \epsilon_i^2$ for a given set of data $(x_{i1}, \dots, x_{ip}, y_i) \in \mathbb{R}^p \times \mathbb{R}$ ($i = 1, \dots, n$), where ϵ_i is the i th residual and defined by $\epsilon_i = y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij}$.

Variable selection in linear regression is the problem to select the best subset of explanatory variables based on a given criterion. In statistical applications, a preferred model keeps the goodness-of-fit for a given data set, and contains as a few unnecessary explanatory variable as possible. In fact, unnecessary explanatory variables may add the noise to the prediction based on the statistical model. As a result, the prediction performance of the model may get worse. In addition, we need to observe and/or monitor more data for unnecessary explanatory variables, and thus will spend more cost due to the unnecessary explanatory variables.

Akaike's information criterion (AIC) is one of criteria for variable selection and proposed in [3]. AIC is used as a measure to select the preferred statistical model in all candidates. The statistical model whose AIC value is the smallest is expected as the preferred statistical mode. In linear regression analysis, this selection corresponds to the selection of a subset of the set of explanatory variables in (1) via AIC. More precisely, for a set $S \subseteq \{1, \dots, p\}$ of candidates of explanatory variables in the statistical model (1), AIC is defined in [3] as follows:

$$\text{AIC}(S) = -2 \max_{\beta, \sigma^2} \{ \ell(\beta, \sigma^2) : \beta_j = 0 \ (j \in \{1, \dots, p\} \setminus S) \} + 2(\#(S) + 2) \quad (2)$$

where $\beta = (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1}$, $\#(S)$ stands for the number of elements in the set S and $\ell(\beta, \sigma^2)$ is the log-likelihood function. Computing AIC values for all subsets S of the explanatory variables in (1), we can obtain the best AIC-based subset. However, since the number of subsets is 2^p , the computation of all subsets is not practical.

Under assumption that all the residual ϵ_i are independent and normally distributed with the zero mean and variance σ^2 , the log-likelihood function can be formulated as

$$\ell(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n \epsilon_i^2.$$

We focus on the first term in (2) to simplify (2). Let S be a set of candidates of explanatory variables in (1). By substituting $\beta_j = 0$ ($j \in \{1, \dots, p\} \setminus S$) to the objective function, the first term can be regarded as the unconstrained minimization. Thus minimum solutions satisfy the following equation

$$\frac{d\ell}{d(\sigma^2)} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n \epsilon_i^2 = 0.$$

From this equation, we obtain $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$. Substituting this equation to (2), we simplify (2) as follows:

$$\begin{aligned} \text{AIC}(S) = \min_{\beta_j} \left\{ n \log \left(\sum_{i=1}^n \epsilon_i^2 \right) : \beta_j = 0 \ (j \in \{1, \dots, p\} \setminus S) \right\} \\ + 2(\#(S) + 2) + n (\log(2\pi/n) + 1). \end{aligned} \quad (3)$$

We use (3) to provide our MINLP formulation of the minimization of AIC in the next section.

The following lemma ensures that the minimization in the first term of (3) has an optimal solution with a finite value.

Lemma 2.1. *For any subset $S \subseteq \{1, \dots, p\}$, the minimization in the first term of (3) has an optimal solution with a finite value.*

Proof. Since the logarithm function has the monotonicity, the optimal solution of the minimization in the first term of (3) is also optimal for the following unconstrained quadratic problem:

$$\min_{\beta_j} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j \in S} \beta_j x_{ij} \right)^2 : \beta_j \in \mathbb{R} (j \in \{0\} \cup S) \right\}. \quad (4)$$

Since the objective function of (4) is bounded below, it follows from [7, Section 9.1.1] that (4) has an optimal solution. \square

3 MINLP formulation for the minimization of AIC

We provide the minimization of $\text{AIC}(S)$ over $S \subseteq \{1, \dots, p\}$ by the following MINLP formulation:

$$\min_{\beta_j, z_j, \epsilon_i, k} \left\{ n \log \left(\sum_{i=1}^n \epsilon_i^2 \right) + 2k : \begin{array}{l} \epsilon_i = y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \ (i = 1, \dots, n), \\ \sum_{j=1}^p z_j = k, \beta_0, \beta_j \in \mathbb{R} \ (j = 1, \dots, p), \\ z_j \in \{0, 1\}, z_j = 0 \Rightarrow \beta_j = 0 \ (j = 1, \dots, p) \end{array} \right\} \quad (5)$$

Here the last constraints represent the logical relationships, *i.e.* β_j has to be zero if $z_j = 0$. This formulation is provided in [12, eq. (22) – (25)].

Next we provide a procedure to find a lower bound of the subproblem of (5) at each node in the branch-and-bound tree. Some variables z_j in (5) are fixed to zero or one at each node of the tree. We define the sets Z_0 , Z_1 and Z for a given node as follows:

$$\begin{aligned} Z_1 &= \{j \in \{1, \dots, p\} : z_j \text{ is fixed to } 1\}, Z_0 = \{j \in \{1, \dots, p\} : z_j \text{ is fixed to } 0\}, \\ Z &= \{j \in \{1, \dots, p\} : z_j \text{ is not fixed}\}. \end{aligned}$$

We remark that $Z_1 \cup Z_0 \cup Z = \{1, \dots, p\}$ and that each set is disjoint with one another. In other words, we can uniquely specify a node in the branch-and-bound search tree by Z_1, Z_0 and Z . We denote the node by $V(Z_1, Z_0, Z)$. Then the subproblem at the node $V(Z_1, Z_0, Z)$ is formulated as follows:

$$\left\{ \begin{array}{l} \min_{\beta_j, z_j} \quad n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) + 2 \sum_{j=1}^p z_j \\ \text{subject to} \quad z_j = 1 \ (j \in Z_1), z_j = 0 \ (j \in Z_0), z_j \in \{0, 1\} \ (j \in Z), \\ \quad \beta_0, \beta_j \in \mathbb{R} \ (j = 1, \dots, p), \beta_j = 0 \ (j \in Z_0) \ z_j = 0 \Rightarrow \beta_j = 0 \ (j \in Z) \end{array} \right. \quad (6)$$

By relaxing the integrality of variables z_j in (6), we obtain the following relaxation problem:

$$\left\{ \begin{array}{l} \min_{\beta_j, z_j} \quad n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) + 2 \sum_{j=1}^p z_j \\ \text{subject to} \quad z_j = 1 \ (j \in Z_1), z_j = 0 \ (j \in Z_0), 0 \leq z_j \leq 1 \ (j \in Z), \\ \beta_0, \beta_j \in \mathbb{R} \ (j = 1, \dots, p), \beta_j = 0 \ (j \in Z_0) \ z_j = 0 \Rightarrow \beta_j = 0 \ (j \in Z) \end{array} \right. \quad (7)$$

Moreover we consider the following problem by eliminating all the logical relationships and all the z_j :

$$\min_{\beta_j} \left\{ n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) + 2\#(Z_1) : \begin{array}{l} \beta_0, \beta_j \in \mathbb{R} \ (j \in Z \cup Z_1), \\ \beta_j = 0 \ (j \in Z_0) \end{array} \right\}. \quad (8)$$

It should be noted that the optimal value of (8) is the same as the optimal value of (7). Hence we deal with (8) as the relaxation problem of (6). In fact, for the optimal solution β^* of (8), we construct a sequence $\{(\beta^N, z^N)\}_{N=1}^\infty$ as follows:

$$\beta^N = \beta^* \text{ and } z_j^N = \begin{cases} 1 & \text{if } j \in Z_1, \\ 1/N & \text{if } j \in Z \text{ and } \beta_j^N \neq 0, \\ 0 & \text{if } j \in Z \text{ and } \beta_j^N = 0, \\ 0 & \text{if } j \in Z_0, \end{cases} \quad (j = 1, \dots, p)$$

for all $N \geq 1$. Clearly, (β^N, z^N) is feasible for (7) for all $N \geq 1$. It is sufficient to prove that the objective value θ^N of (7) at (β^N, z^N) converges the optimal value θ^* of (8) as N goes to ∞ . Since we have

$$\theta^* \leq \theta^N \leq n \log \left(\sum_{i=1}^n \left(y_i - \beta_0^* - \sum_{j=1}^p \beta_j^* x_{ij} \right)^2 \right) + 2\#(Z_1) + \frac{2}{N}\#(Z),$$

the right-hand side converges to θ^* as N goes to ∞ . This implies that the optimal value of (8) is the same as the optimal value of (7).

Although the objective function of (8) contains the logarithm function, we can freely remove the constant $2\#(Z_1)$ and the logarithm by the monotonicity of the logarithm function in (8), and thus obtain the following problem from (8):

$$\min_{\beta_j} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 : \begin{array}{l} \beta_0, \beta_j \in \mathbb{R} \ (j \in Z \cup Z_1) \\ \beta_j = 0 \ (j \in Z_0) \end{array} \right\}. \quad (9)$$

Since (9) is the unconstrained minimization of a quadratic function, we can obtain an optimal solution of (9) by solving a linear system. In our implementation, we call `dposv`, which is a built-in function of LAPACK [4] for solving the linear system. We denote the optimal value of (9) by ξ^* . The optimal value of (8) is $n \log(\xi^*) + 2\#(Z_1)$, which is used as a lower bound of the optimal value of (6).

We provide a procedure that constructs a feasible solution of (5) and computes an upper bound of the optimal value of (5). For this we use an optimal solution $\tilde{\beta} \in \mathbb{R}^{p+1}$ obtained after

solving (9). We define

$$\tilde{z}_j = \begin{cases} 1 & (\text{if } j \in \tilde{Z} \cup Z_1), \\ 0 & (\text{otherwise}) \end{cases} \quad (j = 1, \dots, p), \tilde{\epsilon}_i = y_i - \tilde{\beta}_0 - \sum_{j=1}^p \tilde{\beta}_j x_{ij} \quad (i = 1, \dots, n) \text{ and } \tilde{k} = \sum_{j=1}^p \tilde{z}_j,$$

where $\tilde{Z} = \{j \in Z : \tilde{\beta}_j \neq 0\}$. It is easy to see that $(\tilde{\beta}_j, \tilde{z}_j, \tilde{\epsilon}_i, \tilde{k})$ is feasible for (5) and the objective value is $n \log(\xi^*) + 2\#(\tilde{Z} \cup Z_1)$. If the objective value is smaller than the current best upper bound, then we update the current best upper bound.

Finally, we give another understanding for our proposed formulation and propose an efficient computation based on this understanding.

- Since we can regard (9) as linear regression whose explanatory variables are in $Z_1 \cup Z$, the computation of the lower bound from (9) corresponds to the computation of the value $\text{AIC}(Z_1 \cup Z) - 2\#(Z)$, while the upper bound corresponds to the AIC value of the statistical model whose explanatory variables are in $Z_1 \cup Z$, *i.e.* $\text{AIC}(Z_1 \cup Z)$. Therefore, our proposed method computes the AIC value of the the statistical model with $Z_1 \cup Z$ at each node $V(Z_1, Z_0, Z)$, up to constant term $4 + n(\log(2n\pi) + 1)$ of (3). In summary, we consider that our proposed method branches and prunes the branch-and-bound search tree efficiently by using this understanding.
- The statistical package `leaps` [11] in R [14] adopts the branch-and-bound scheme in a similar manner. A QR decomposition is exploited at each node in the branch-and-bound search tree. In particular, `leaps` solve a linear system effectively by using the QR decomposition obtained at its parent node.

`leaps` finds the best statistical model much faster than our proposed method for data sets whose p is less than or equal to 32 and which do not have linear dependency introduced in subsection 4.2. If the data set has linear dependency, `leaps` does not work effectively, while our proposed method works more efficiently by using the linear dependency in data sets. This technique will be discussed in Section 4.2.

- We provide an efficient computation of lower and upper bounds based on this understanding. We assume that we obtain the lower and upper bounds at a node $V(Z_1, Z_0, Z)$. Then we do not need to solve (9) at its child node $V(Z_1 \cup \{j\}, Z_0, Z \setminus \{j\})$, where $j \in Z$. This node is generated by branching $z_j = 1$ at the node $V(Z_1, Z_0, Z)$. In fact, since we have $(Z_1 \cup \{j\}) \cup (Z \setminus \{j\}) = Z_1 \cup Z$, the relaxation problem (9) at the child node $V(Z_1 \cup \{j\}, Z_0, Z \setminus \{j\})$ is equivalent to one at the node $V(Z_1, Z_0, Z)$. Thus the upper bound at the child node is the same as one at the node $V(Z_1, Z_0, Z)$, and the lower bound is the lower bound computed at the node $V(Z_1, Z_0, Z)$ plus two because of $2\#(Z_1 \cup \{j\}) = 2\#(Z_1) + 2$.

4 Some techniques to improve the numerical performance

We describe some techniques to improve numerical performance to solve (5).

4.1 SCIP

In order to implement our proposed method, we use SCIP [2, 13, 17], which is a mathematical optimization software and a branch-and-bound framework. In fact, it has high user plug-in flexibility which helps to solve (5) efficiently. We implement a procedure, which is called

relaxator or *relaxation handler*, to obtain lower bounds as in Section 3. In addition, we also implement procedures to compute upper bounds via a method based on stepwise methods discussed in subsection 4.3 and to define branching rules described in subsection 4.4.

4.2 Handling the linear dependency in data

We illustrate that we can efficiently compute the optimal value of (5) by using the linear dependency in data. Although linearly independent data is often the assumption in standard statistical textbooks, practical data has often linear dependency, *e.g.* `servo` and `auto-mpg` in UCI Machine Learning Repository [16].

For a set of given data $(x_{i1}, \dots, x_{ip}, y_i) \in \mathbb{R}^p \times \mathbb{R}$ ($i = 1, \dots, n$), we denote

$$x^0 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, x^j = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{pmatrix} \quad (j = 1, \dots, p).$$

We say that data has linear dependent variables if the vectors $x^0, x^1, \dots, x^p \in \mathbb{R}^n$ are linearly dependent.

From the definition of the linear dependency in data, we can reduce the computational cost for solving (9) when the data has linearly dependency. At a node $V(Z_1, Z_0, Z)$, if there exists a subset $S \subseteq Z_1 \cup Z$ such that the vectors $\{x^k : k \in S \cup \{0\}\}$, we can fix one of variables z_j in $j \in S \cap Z$ to zero. In fact, since we have $\sum_{j \in S \cup \{0\}} \alpha_j x^j = 0$ for some $(\alpha_j)_{j \in S \cup \{0\}} \neq 0$, we can remove one variable z_j by substituting this equation to (9). This implies that the number of variables in (9) decrease, and thus we solve the linear equation with a fewer variables.

Moreover we can prune some nodes efficiently by using the linear dependency. The following lemma ensures that we do not need to branch $z_q = 1$ for some $q \in Z$ if the data has the linear dependency. Thus we need to handle only $z_q = 0$ in this case.

Lemma 4.1. *Assume that in (6), there exists $q \in Z$ such that the vector x^q and vectors $\{x^j : j \in Z_1 \cup \{0\}\}$ are linearly dependent. Then an optimal solution of (6) satisfies $z_q = 0$.*

Proof. Let $(\tilde{\beta}_j, \tilde{z}_j)$ be an optimal solution of (6), and θ^* be the optimal value of (6). Suppose that $\tilde{z}_q = 1$. It follows from the assumption that there exists $\alpha_j \in \mathbb{R}$ ($j \in Z_1 \cup \{0\}$) such that $(\alpha_j)_{j \in Z_1 \cup \{0\}} \neq 0$ and

$$x^q = \sum_{j \in Z_1 \cup \{0\}} \alpha_j x^j.$$

Then the following solution $(\hat{\beta}_j, \hat{z}_j)$ is also feasible for (6):

$$\hat{\beta}_j = \begin{cases} \tilde{\beta}_j + \tilde{\beta}_q \alpha_j & (\text{if } j \in (Z \setminus \{q\}) \cup Z_1 \cup \{0\}), \\ 0 & (\text{otherwise}) \end{cases}, \quad \text{and } \hat{z}_j = \begin{cases} 1 & (\text{if } j \neq q \text{ and } \tilde{z}_j = 1), \\ 0 & (\text{otherwise}) \end{cases}$$

The objective value of (6) at $(\hat{\beta}_j, \hat{z}_j)$ is $\theta^* - 2$, which contradicts the optimal value θ^* . \square

A given set of data which has linear dependency satisfies the assumption of Lemma 4.1. In fact, there exists a subset $S \subseteq \{1, \dots, p\}$ such that the vectors $\{x^k : k \in S \cup \{0\}\}$ are linearly dependent. Hence Lemma 4.1 ensures that we do not need to generate a child node by branching $z_q = 1$ at a node $V(Z_1, Z_0, Z)$ when $q \in S \cap Z$ and $S \setminus \{q\} \subseteq Z_1$.

In addition, if there exists a subset $S \subseteq \{1, \dots, p\}$ such that for every $j \in S$, the vectors $\{x^k : k \in \{0\} \cup (S \setminus \{j\})\}$ are linearly dependent, then we can prune some nodes before applying our proposed method to (5). In fact, it follows from the assumption on S that for every $j \in S$ we do not need to branch $z_j = 1$ at the node $V(Z_1, Z_0, Z)$. This implies that optimal solutions of (5) satisfy the following linear inequality:

$$\sum_{j \in S} z_j \leq \#(S) - 1.$$

By adding this inequality in (5), we do not generate any nodes in which $S \subseteq Z_1$ hold. We execute a greedy algorithm in Algorithm 1 to find a collection \mathcal{C} of such sets S .

Algorithm 1: A greedy algorithm to find a collection of sets of linearly dependent vectors

Input: Data $x^0, x^1, x^2, \dots, x^p \in \mathbb{R}^n$

Output: A collection \mathcal{C} of sets of linearly dependent vectors

$\mathcal{C} \leftarrow \emptyset, S \leftarrow \emptyset;$

for $j \rightarrow 0$ **to** p **do**

if the vectors $\{x^j : j \in \{0\} \cup S \cup \{j\}\}$ is linearly independent **then**

$S \leftarrow S \cup \{j\};$

else

 Solve the following linear equation:

$$\sum_{k \in S \cup \{0\}} \alpha_k x^k = x^j. \quad (10)$$

$S' \leftarrow \{k \in S : \alpha_k \neq 0\}, \mathcal{C} \leftarrow \mathcal{C} \cup \{S'\};$

end

end

return $\mathcal{C};$

We remark that the linear equation (10) has a unique solution because the matrix $(x^k)_{k \in S \cup \{0\}}$ is of full column rank.

4.3 Computation of upper bounds based on stepwise methods

Although we mainly use the procedure described in Section 3 to compute upper bounds, we also use the stepwise methods with forward selection (SW_+) and backward elimination (SW_-). SW_+ starts with no explanatory variables and adds one explanatory variable at a time until the AIC value does not decrease. More precisely, for the current set S of explanatory variables, we choose an explanatory variable whose the AIC value $\text{AIC}(S \cup \{j\})$ is minimized over $j \in \{1, \dots, p\} \setminus S$. SW_- is just the reverse of SW_+ . It starts with all explanatory variables and remove one explanatory variable at a time until the AIC value does not decrease. Note that since these methods add or remove one explanatory variable at a time, they may miss the best statistical model. In this sense, we can say that they are local search algorithms for variable selection.

We describe our heuristics to computer an upper bound in more details in Algorithm 2. To

this end, we define $S \subseteq Z_1 \cup Z$ for subproblem (6) and consider the following problem:

$$\begin{cases} \min_{\beta_j, z_j} & n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) + 2 \sum_{j=1}^p z_j \\ \text{subject to} & \beta_0, \beta_j \in \mathbb{R}, z_j = 1 \ (j \in S), \beta_j = 0, z_j = 0 \ (j \in \{1, \dots, p\} \setminus S) \end{cases} \quad (11)$$

We denote the optimal value and an optimal solution of (11) by $\bar{\theta}_S$ and $(\bar{\beta}_S, \bar{z}_S)$, respectively.

Algorithm 2: Stepwise methods to compute an upper bound

Input: Z_1, Z_0 and Z
Output: A feasible solution (β, z) of (6)
 /* Stepwise method with forward selection */
 $S \leftarrow Z_1, v_f \leftarrow \infty;$
while $\bar{\theta}_S < v_f$ **do**
 | $v_f \leftarrow \bar{\theta}_S, (\beta_f, z_f) \leftarrow (\bar{\beta}_S, \bar{z}_S);$
 | Find $j \in Z \setminus S$ such that $\theta_{S \cup \{j\}}$ is minimized over all $j \in Z \setminus S;$
 | $S \leftarrow S \cup \{j\};$
end
 /* Stepwise method with backward elimination */
 $S \leftarrow Z_1 \cup Z, v_b \leftarrow \infty;$
while $\bar{\theta}_S < v_b$ **do**
 | $v_b \leftarrow \bar{\theta}_S, (\beta_b, z_b) \leftarrow (\bar{\beta}_S, \bar{z}_S);$
 | Find $j \in Z \cap S$ such that $\theta_{S \setminus \{j\}}$ is minimized over all $j \in Z \cap S;$
 | $S \leftarrow S \setminus \{j\};$
end
if $v_f < v_b$ **then**
 | **return** $(\beta_f, z_f);$
else
 | **return** $(\beta_b, z_b);$
end

We remark that an optimal solution of (11) is feasible for the subproblem (6) if $Z_1 \subseteq S$. Since S always contains Z_1 in Algorithm 2, the returned solution (β, z) is feasible for (6). In addition, we set Z_1 as the initial set of SW_+ instead of the empty set because we execute Algorithm 2 at the node $V(Z_1, Z_0, Z)$. Similarly, we set $Z_1 \cup Z$ as the initial set of SW_- . These are different from the original stepwise methods.

In statistical applications, instead of finding the global minimum of (5), stepwise methods, which are local search algorithms, are commonly used in practice. In fact, they often find a better statistical model and work effectively in our implementation. However since stepwise methods spend more computational costs than the procedure described in Section 3, we apply Algorithm 2 to only subproblem (6) at the node whose depth from the root node is less than or equal to 10 in our implementation.

4.4 Most frequent branching and Strong branching

We define two branching rules for variables z_j to improve the performance of our implementation. The first one is called *most frequent branching* and uses all stored feasible solutions in the procedure to compute upper bounds. The second one is called *strong branching*. This is based

on the strong branching rule in [1, Section 5.4]. We propose a more efficient computation for the strong branching rule than [1]. We will show the numerical comparison with branching rules implemented in SCIP in subsection 5.2. We will observe from the numerical results that most frequent branching is effective for a set data which has linear dependency, while strong branching is effective for a set data which does not have linear dependency.

The most frequent branching is based on the tendency that some explanatory variables adopted for the best statistical model are also used in statistical models whose AIC value is close to the smallest AIC value. By branching variables z_j in (6) which correspond to such explanatory variables, we can expect that (6) at the node generated by $z_j = 0$ is pruned as early as possible. To find such explanatory variables, we use feasible solution stored in our procedure to compute upper bounds. We describe the most frequent branching rule at the current node in Algorithm 3.

Algorithm 3: Most frequent branching rule

Input: A positive integer N , a set Z of unfixed variables in the node and all feasible solutions of (5) found from the root node through the current node

Output: $J \in Z$

Choose N feasible solutions $(\beta^1, z^1), \dots, (\beta^N, z^N)$ out of all stored feasible solutions;

/ Here (β^i, z^i) is a feasible solution of (5) whose objective value is the i th smallest in all the stored solutions */*

for $j \in Z$ **do**

 | Compute score value s_j defined by $s_j = \#(T_j)$, where $T_j = \{\ell \in \{1, \dots, N\} : z_j^\ell = 1\}$;

end

return $J \in Z$ with $s_J = \max_{j \in Z} \{s_j\}$;

We observe in our preliminary numerical experiment that the obtained lower bound at the child node generated by $z_j = 0$ tends to be relatively bigger and that the pruning process tends to work earlier in comparison to branching rules of SCIP. As a result, our proposed method with the most frequent branching rule often visits a fewer nodes in the branch-and-bound tree.

In the strong branching rule, we compute lower bounds for all possible branching $z_k = 1$ and $z_k = 0$, and choose index $k \in Z$ so that the lower bound is maximized in all computed lower bounds. More precisely, for the subproblem (6) at a node $V(Z_1, Z_0, Z)$ and $k \in Z$, the relaxation problem of the subproblem branched by $z_k = 1$ and $z_k = 0$ can be formulated as (12) and (13) as follows, respectively.

$$\begin{cases} \min_{\beta_j} & n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) + 2\#(Z_1 \cup \{k\}) \\ \text{subject to} & \beta_0, \beta_j \in \mathbb{R} \ (j \in (Z \setminus \{k\}) \cup (Z_1 \cup \{k\})), \beta_j = 0 \ (j \in Z_0) \end{cases} \quad (12)$$

$$\begin{cases} \min_{\beta_j} & n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) + 2\#(Z_1) \\ \text{subject to} & \beta_0, \beta_j \in \mathbb{R} \ (j \in (Z \setminus \{k\}) \cup Z_1), \beta_j = 0 \ (j \in Z_0 \cup \{k\}) \end{cases} \quad (13)$$

Since we have $(Z \setminus \{k\}) \cup (Z_1 \cup \{k\}) = Z \cup Z_1$, the optimal value of (12) for all $k \in Z$ is $\theta^* + 2$, where θ^* is the optimal value of (8) at a node $V(Z_1, Z_0, Z)$. Hence we select an index $k \in Z$ only from all optimal values θ_k^* of (13). We describe the strong branching rule at the current node in Algorithm 4.

Algorithm 4: Strong branching rule

Input: Subproblem (6) in the node $V(Z_1, Z_0, Z)$
Output: $J \in Z$
for $k \in Z$ **do**
 | Solve (13) with k and obtain optimal value θ_k^* ;
end
return $J \in Z$ with $\theta_J^* = \max_{k \in Z} \{\theta_k^*\}$;

5 Numerical experiments

We implement our approach and procedures discussed in Sections 3 and 4, and apply our implementation¹ to benchmark data sets in [16]. We apply our implementation to standardized data sets, *i.e.* the data is transformed to have the zero mean and unit variance. Note that the standardized data has also linear dependency even if we apply the standardization to the original data which has linear dependency. The specification of the computer is CPU : 3.5 GHz Intel Core i7, Memory : 16GB and OS : OS X 10.9.5. In subsections 5.1 and 5.3, we adopt the most frequent branching rule for data which has linear dependency, while we adopt the strong branching for data which does not have linear dependency. In subsection 5.2, we discuss the reason why we use the different branching rules.

5.1 Comparison with stepwise methods and MISOCP approach

We compare our proposed method with stepwise methods (SW_+ and SW_-) and the MISOCP approach proposed in [12] via CPLEX [9]. This approach is also obtained from (5). Although the objective function of (5) is non-convex, the difficulty due to the non-convexity is overcome by using the identity $\exp(\log(x)) = x$ and the monotonicity of the exponential function $\exp(x)$. See [12, Section 3.2] for the detail. The resulting problem is formulated as MISOCP and is tractable by CPLEX.

Table 1 shows the summary of numerical comparisons. The mark \bullet in the first column indicates that the data has linear dependency. The second, third, and sixth columns indicate the numbers of data, the explanatory variables in the statistical model (1), and the ones in the models found by using each method. The fifth column indicates the obtained AIC values by each method. The values with the bold font are the best among four values. The seventh column indicates the cpu time in seconds to compute the optimal value. “>5000” means that the corresponded method cannot find the optimal value within 5000 seconds. The last column indicates the gap in the percent as follows:

$$\text{gap} = \frac{\text{upper bound} - \text{lower bound}}{\max\{1, |\text{upper bound}|\}} \times 100.$$

It should be noted that if the gap is sufficiently close to zero, then the obtained value is optimal. MINLP, MISOCP, SW_+ and SW_- indicate the results obtained by our proposed method, MISOCP approach and the stepwise method with forward selection and backward elimination, respectively. We observe the following from Table 1.

¹This is available at [18].

- MINLP computes the optimal value much faster than MISOCP. MINLP finds smaller AIC values than MISOCP even when MINLP cannot find them within 5000 seconds.
- The AIC value obtained by SW_+ or SW_- is equal to one by MINLP, *i.e.* **crime** and **forestfires**. In fact, as we mentioned in subsection 4.3, we use stepwise methods in some nodes in our implementation. This implies that our procedure to compute an upper bound discussed in Section 3 cannot find better feasible solutions than ones by the stepwise methods.

5.2 Comparison of branching rules

We compare the numerical performance of the most frequent branching and strong branching with branching rules implemented in SCIP. In Table 2, Std, MFB and SB stand for numerical results by the branching rules in SCIP, the most frequent branching rule and the strong branching rule. The sixth column indicates the number of visited nodes by our proposed method with the applied branching rule. The values with the bold font are the best among three values. We observe from Table 2:

- The most frequent branching rule works more effectively than other ones for sets of data which have linear dependency. In fact, the gap by the most frequent branching rule is the smallest and the computation time is the shortest. In addition, The number of the visited nodes by the most frequent branching is also smaller than other branching rules. In contrast, the strong branch is more efficient than other branching rules for data which do not have linear dependency.
- For $p \leq 32$, the gap obtained by the best branching rule is the smallest in three branching rules, though it visits fewest nodes in the branch-and-bound tree. This means that the best branching computes tighter lower bounds than other branching rules.
- These are the reasons why we use different branching rules in Tables 1 and 3.

5.3 Comparison with MIQP formulation

Bertsimas and Shioda [6] and Bertsimas et al [5] provide a mixed integer quadratic programming (MIQP) formulation with a cardinality constraint for linear regression. Their formulation is available to the minimization of AIC by fixing the number of explanatory numbers from 0 to p . In fact, the minimization can be equivalently reformulated as follows:

$$\min_{k=0,\dots,p} \min_{S \subseteq \{1,\dots,p\}} \{AIC(S) : \#(S) = k\}. \quad (14)$$

Since each inner optimization problem in (14) can be formulated as a MIQP problem, we can obtain the best statistical model by solving all $(p + 1)$ optimization problems. In this subsection, we introduce a MIQP formulation by Bertsimas and Shioda [6] and Bertsimas et al [5] for the inner optimization problems in (14). In addition, we provide a more efficient algorithm than this naive algorithm and compare the algorithm with our proposed method.

Name	n	p	Methods	AIC	k	time(sec)	gap(%)
housing	506	13	MINLP	776.21	11	0.04	0.00
			MISOCP	776.21	11	7.96	0.00
			SW ₊	776.21	11	0.35	–
			SW _–	776.21	11	0.10	–
●servo	167	19	MINLP	258.35	9	0.79	0.00
			MISOCP	258.35	9	7.99	0.00
			SW ₊	258.35	9	0.19	–
			SW _–	260.16	10	0.18	–
●auto-mpg	392	25	MINLP	332.88	15	1.76	0.00
			MISOCP	332.88	15	303.83	0.00
			SW ₊	334.73	16	0.49	–
			SW _–	337.96	18	0.32	–
●solarflareC	1066	26	MINLP	2816.29	9	10.49	0.00
			MISOCP	2816.29	9	304.51	0.00
			SW ₊	2816.29	9	0.45	–
			SW _–	2821.61	12	1.08	–
●solarflareM	1066	26	MINLP	2926.90	7	3.99	0.00
			MISOCP	2926.90	7	255.02	0.00
			SW ₊	2926.90	7	0.36	–
			SW _–	2930.91	9	1.16	–
●solarflareX	1066	26	MINLP	2882.80	3	0.92	0.00
			MISOCP	2882.80	3	19.39	0.00
			SW ₊	2882.80	3	0.18	–
			SW _–	2891.56	9	1.20	–
breastcancer	194	32	MINLP	508.40	10	90.21	0.00
			MISOCP	508.62	10	>5000	3.72
			SW ₊	509.50	8	0.24	–
			SW _–	509.96	14	0.60	–
●forestfires	517	63	MINLP	1429.64	12	>5000	0.77
			MISOCP	1431.32	12	>5000	6.44
			SW ₊	1429.64	12	0.94	–
			SW _–	1447.36	21	7.43	–
●automobile	159	65	MINLP	-61.28	32	>5000	13.95
			MISOCP	-55.83	34	>5000	27.22
			SW ₊	-28.55	21	1.12	–
			SW _–	-47.61	40	2.64	–
crime	1993	100	MINLP	3410.25	50	>5000	0.50
			MISOCP	3469.34	74	>5000	8.51
			SW ₊	3430.19	37	17.03	–
			SW _–	3410.25	50	105.40	–

Table 1: Summary of numerical results by MINLP, MISOCP, SW₊ and SW_–

Name	Methods	AIC	k	time(sec)	# of visited nodes	gap(%)
housing	Std	776.21	11	0.05	55	0.00
	MFB	776.21	11	0.05	49	0.00
	SB	776.21	11	0.04	27	0.00
•servo	Std	258.35	9	1.17	7577	0.00
	MFB	258.35	9	0.79	4705	0.00
	SB	258.35	9	0.41	2261	0.00
•auto-mpg	Std	332.88	15	4.06	18959	0.00
	MFB	332.88	15	1.76	5723	0.00
	SB	332.88	15	2.68	11586	0.00
•solarflareC	Std	2816.29	9	53.33	166639	0.00
	MFB	2816.29	9	10.49	32261	0.00
	SB	2816.29	9	23.13	79015	0.00
•solarflareM	Std	2926.90	7	40.03	117889	0.00
	MFB	2926.90	7	3.99	11903	0.00
	SB	2926.90	7	23.72	81899	0.00
•solarflareX	Std	2882.80	3	4.37	9737	0.00
	MFB	2882.80	3	0.92	1519	0.00
	SB	2882.80	3	3.40	7453	0.00
breastcancer	Std	508.40	10	505.70	3851×10^3	0.00
	MFB	508.40	10	478.66	3422×10^3	0.00
	SB	508.40	10	90.21	550×10^3	0.00
•forestfires	Std	1429.64	12	>5000	7480×10^3	1.11
	MFB	1429.64	12	>5000	13179×10^3	0.77
	SB	1429.64	12	>5000	9938×10^3	0.95
•automobile	Std	-60.29	32	>5000	32192×10^3	12.30
	MFB	-61.28	32	>5000	29785×10^3	13.95
	SB	-61.59	33	>5000	15300×10^3	16.43
crime	Std	3410.25	50	> 5000	10272×10^3	0.78
	MFB	3410.25	50	> 5000	9753×10^3	0.52
	SB	3410.25	50	> 5000	1904×10^3	0.50

Table 2: Summary of numerical results by branching rules in SCIP (Std), the most frequent branching (MFB) and strong branching (SB)

Each inner optimization problem in (14) can be reformulated as follows:

$$\min_{\beta_j} \left\{ \begin{array}{l} \epsilon_i = y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \quad (i = 1, \dots, n), \\ n \log \left(\sum_{i=1}^n \epsilon_i^2 \right) + 2k : \sum_{j=1}^p z_j = k, \beta_0, \beta_j \in \mathbb{R} \quad (j = 1, \dots, p), \\ z_j \in \{0, 1\}, z_j = 0 \Rightarrow \beta_j = 0 \quad (j = 1, \dots, p) \end{array} \right\} \quad (15)$$

For any fixed k , since the logarithm function in (15) has the monotonicity, we can find an optimal solution (15) by solving the following quadratic programming problem:

$$\min_{\beta_j} \left\{ \begin{array}{l} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 : \sum_{j=1}^p z_j = k, \beta_0 \in \mathbb{R}, \\ z_j \in \{0, 1\} \quad (j = 1, \dots, p), \\ z_j = 0 \Rightarrow \beta_j = 0 \quad (j = 1, \dots, p) \end{array} \right\} \quad (16)$$

(16) is a MIQP formulation. We denote the optimal value of (16) by η_k^* . If (16) is infeasible, we set $\eta_k^* = +\infty$. Then the optimal value of inner problem (15) with k is $n \log(\eta_k^*) + 2k$. Therefore we obtain the optimal value and solution of (14) by computing all optimal values of (15) for $k = 0, \dots, p$. We describe the naive algorithm in Algorithm 5.

Algorithm 5: Naive algorithm for (5) via MIQP

Input: Minimization of AIC (5)

Output: An optimal solution of (5)

for $k \rightarrow 0$ **to** p **do**

 | Find the optimal value η_k^* and an optimal solution (β_k^*, z_k^*) of (16) with k ;

end

Find an index K with $\theta_K^* = \min_{k=0, \dots, p} \{n \log(\eta_k^*) + 2k\}$;

return (β_K^*, z_K^*) ;

The following lemma ensures that we can find an upper bound of k if we have a feasible solution of (3).

Lemma 5.1. *Let $\hat{\theta} \in \mathbb{R}^{p+1}$ be the optimal value of the following optimization problem:*

$$\min_{\beta_j} \left\{ n \log \left(\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right) : \beta_0, \beta_j \in \mathbb{R} \quad (j = 1, \dots, p) \right\}. \quad (17)$$

In addition, $\bar{\theta}$ is the objective value of (5) at a feasible solution of (5). Then any optimal solution (β^, z^*) of (5) satisfies*

$$\sum_{j=1}^p z_j^* \leq \left\lfloor \frac{\bar{\theta} - \hat{\theta}}{2} \right\rfloor.$$

Proof. Let θ^* be the optimal value of (5) and (β^*, z^*) be an optimal solution of (5). Then we have

$$\bar{\theta} \geq \theta^* = n \log \left(\sum_{i=1}^n \left(y_i - \beta_0^* - \sum_{j=1}^p \beta_j^* x_{ij} \right)^2 \right) + 2 \sum_{j=1}^p z_j^* \geq \hat{\theta} + 2 \sum_{j=1}^p z_j^*,$$

and thus we have $\sum_{j=1}^p z_j^* \leq (\bar{\theta} - \hat{\theta})/2$. Since z_j^* is integer, we obtain the desired result. \square

We describe an algorithm based on Lemma 5.1 in Algorithm 6.

Algorithm 6: Faster algorithm for (5) via MIQP

Input: Minimization of AIC (5)
Output: An optimal solution of (5)
Solve (17) and obtain $\hat{\theta}$;
 $\bar{\theta} \leftarrow +\infty$;
for $k \rightarrow 0$ **to** p **do**
 if $k > \lfloor \frac{\bar{\theta} - \hat{\theta}}{2} \rfloor$ **then**
 | Stop;
 end
 Find the optimal value η_k^* and solution (β_k^*, z_k^*) of (16) with k ;
 if $\bar{\theta} \geq n \log(\eta_k^*) + 2k$ **then**
 | $\bar{\theta} \leftarrow n \log(\eta_k^*) + 2k, (\beta^*, z^*) \leftarrow (\beta_k^*, z_k^*)$;
 end
end
return (β^*, z^*) ;

We give details on our numerical experiment.

- We solve (15) by CPLEX. In particular, since the last constraints in (15) represent the logical relationship between z_j and β_j , we use *indicator* implemented in CPLEX to represent these constraints.
- We add linear inequalities in (15) by applying Lemma 4.1 to (15) when a given set of data has linear dependency. See subsection 4.2 for the detail.
- We also solve optimization problems obtained by replacing the constraint $\sum_{j=1}^p z_j = k$ by $\sum_{j=1}^p z_j \leq k$ in (15). In Table 3, “Fast \leq ” indicates that we solve those problems in Algorithm 6, while “Fast=” indicates that we solve (15) in Algorithm 6. By this replacement, we can use an optimal solution (β_k^*, z_k^*) of the optimization problem with k to compute an upper bound of the optimization problem with $k + 1$.
- We terminate if the corresponded method cannot find the best AIC value within 5000 seconds. In addition, the values with the bold font are the best among four values except for “>5000” in the last column.

We provide numerical results on our proposed method, Algorithms 5 and 6 in Table 3. We observe the following from Table 3:

- MINLP outperforms MIQP approaches. In particular, for larger p , MINLP obtains much better AIC values than MIQP approaches although all approaches cannot solve within 5000 seconds.
- The performance of Fast \leq is similar to Fast=, though Fast \leq uses an initial upper bound.

6 Conclusion

We propose the MINLP formulation (5) of AIC minimization for linear regression, and implement it by using SCIP. We formulate an unconstrained optimization problem (8) as the

Name	Methods	AIC	k	time(sec)
housing	MINLP	776.21	11	0.04
	Naive	776.21	11	2.54
	Fast=	776.21	11	2.15
	Fast \leq	776.21	11	2.43
•servo	MINLP	258.35	9	0.79
	Naive	258.35	9	2.27
	Fast=	258.35	9	1.27
	Fast \leq	258.35	9	1.29
•auto-mpg	MINLP	332.88	15	1.76
	Naive	332.88	15	22.22
	Fast=	332.88	15	19.04
	Fast \leq	332.88	15	14.45
•solarflareC	MINLP	2816.29	9	10.49
	Naive	2816.29	9	26.49
	Fast=	2816.29	9	18.17
	Fast \leq	2816.29	9	15.03
•solarflareM	MINLP	2926.90	7	3.99
	Naive	2926.90	7	25.27
	Fast=	2926.90	7	8.15
	Fast \leq	2926.90	7	7.24
•solarflareX	MINLP	2882.80	3	0.92
	Naive	2882.80	3	10.65
	Fast=	2882.80	3	2.25
	Fast \leq	2882.80	3	2.40
breastcancer	MINLP	508.40	10	90.21
	Naive	508.40	10	420.44
	Fast=	508.40	10	402.64
	Fast \leq	508.40	10	421.96
•forestfires	MINLP	1429.64	12	>5000
	Naive	1435.07	7	>5000
	Fast=	1435.07	7	>5000
	Fast \leq	1435.07	7	>5000
•automobile	MINLP	-61.28	32	>5000
	Naive	52.84	8	>5000
	Fast=	52.84	8	>5000
	Fast \leq	52.84	8	>5000
crime	MINLP	3410.25	50	>5000
	Naive	3646.35	4	>5000
	Fast=	3646.35	4	>5000
	Fast \leq	3646.35	4	>5000

Table 3: Summary of numerical results by MINLP, Naive (Algorithm 5), Fast= and Fast \leq (Algorithm 6)

relaxation problem of the subproblem (6). As a result, a lower bound can be computed by solving a linear equation at each node. In addition, an upper bound is the lower bound plus a constant, and a feasible solution is generated from a solution after solving the relaxation problem (8).

We implement this procedure with SCIP because it has the high flexibility in the user plugin. In fact, we implement a relaxator to compute lower and upper bounds, and two branching rules to prune subproblems efficiently. In addition, our implementation efficiently prunes and branches subproblems by using linear dependency in data set and two branching rules. As a result, we can obtain the best statistical models (1) for $p \leq 32$. In addition, we observe that our implementation outperforms MISOCP approach [12] and MIQP approaches [6, 5] in our numerical experiments.

Future work involves to apply our implementation to data sets with larger p and/or n . A possible choice to accomplish this involves the use of parallel computation via ParaSCIP and FiberSCIP [15]. Secondly, various non-AIC information criterion, *e.g.* BIC and Hannan-Quinn information criteria are already proposed. By changing the objective function in (5), our proposed method can be applied to these information criteria as well.

Acknowledgements

The second author was supported by JSPS KAKENHI Grant Numbers 26400203.

References

- [1] T. Achterberg, “Constraint Integer Programming”, Ph.D. Thesis, Technische Universität Berlin, 2007.
- [2] T. Achterberg, “SCIP: solving constraint integer programs”, *Math. Prog. Comp.*, 1, 1, 1 – 41, 2009.
- [3] H. Akaike, “A new look at the statistical model identification”, *IEEE Trans. Autom. Control*, 19, 6, 716–723, 1974.
- [4] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen, “LAPACK Users’ Guide”, SIAM, 1999.
- [5] D. Bertsimas, A. King and R. Mazumder, “Best Subset Selection via a Modern Optimization Lens”, *Ann. Stat.*, 44, 2, 813 – 852, 2016.
- [6] D. Bertsimas and R. Shioda, “Algorithm for cardinality-constrained quadratic optimization”, *Comput. Optim. Appl.*, 43, 1 – 22, 2009.
- [7] S. Boyd and L. Vandenberghe, “Convex Optimization”, Cambridge University Press, 2004.
- [8] I. Guyon and A. Elisseeff, “An Introduction to Variable and Feature Selection”, *J. Mach. Learn. Res.*, 3, 1157 – 1182, 2003.
- [9] IBM ILOG CPLEX Optimizer 12.6.2, IBM ILOG 2015.

- [10] K. Kimura and H. Waki, “Minimization of Akaike’s Information Criterion via Mixed Integer Nonlinear Program”, to appear in Proceedings in the 5th International Congress on Mathematical Software, Berlin, 2016.
- [11] T. Lumley, “Package ‘leaps’”, 2015.
- [12] R. Miyashiro and Y. Takano, “Mixed integer second-order cone programming formulations for variable selection”, *Eur. J. Oper. Res.*, 247, 721 – 731, 2015.
- [13] SCIP: Solving Constraint Integer Programs, <http://scip.zib.de/>
- [14] R Development Core Team, “R: A Language and Environment for Statistical Computing”, R Foundation for Statistical Computing, Vienna, Austria, 2008, <http://www.R-project.org>
- [15] Y. Shinano, T. Achterberg, T. Berthold, S. Heinz and T. Koch, “ParaSCIP – a parallel extension of SCIP”, *Competence in High Performance Computing 2010*, editors: C. Bischof, H.-G. Hegering, W. E. Nagel and G. Wittum, 135 – 148, Springer, 2012.
- [16] UCI Machine Learning Repository, <http://archive.ics.uci.edu/ml/>
- [17] S. Vigerske and A. Gleixner, “SCIP: Global Optimization of Mixed-Integer Nonlinear Programs in a Branch-and-Cut Framework”, ZIB-Report 16-24, Zuse Institute Berlin, May 2016.
- [18] <https://github.com/k-kimura1224/MAIC>