

# Variable Selection with Big Data based on Zero Norm and via Sequential Monte Carlo

Jin-Chuan Duan\*

(First Version: October 23, 2018; This Version: April 22, 2019)

## Abstract

Selecting a subset from many potential explanatory variables in linear regressions has long been the subject of research interest, and the matter is made more important in the era of big data when many more variables become available/accessible. Of late, the  $l_1$ -norm penalty based techniques such as Lasso of Tibshirani (1996) have become very popular. However, the variable selection problem in its natural setting is a zero-norm penalty problem, i.e., a penalty on the number of variables as opposed to the  $l_1$ -norm of the regression coefficients. The popularity of the  $l_1$ -norm penalty or its variants has more to do with computational considerations, because selection with the zero-norm penalty is a highly demanding combinatory optimization problem when the number of potential variables becomes large. We devise a sequential Monte Carlo (SMC) method as a practical and reliable tool for zero-norm variable selection problems, and selecting, say, best 20 out of 1,000 potential variables can, for example, be completed with a typical multi-core desktop computer in a couple of minutes. The methodological essence is to understand that the selection problem is equivalent to the task of sampling from a discrete probability function defined over all possible combinations comprising, say,  $k$  regressors out of  $p \geq k$  potential variables, where the peak of this function corresponds to the optimal combination. The solution technique sets out to sequentially generate samples, and after a while the final sample represents the target probability function. With the final SMC sample in place, we deploy the extreme value theory to assess how likely and to what extent the maximum  $R^2$  has been achieved. We also demonstrate through a simulation study the method's reliability and superiority vis-a-vis the adaptive Lasso.

---

\*Duan is with the National University of Singapore (Business School, Risk Management Institute and Department of Economics). E-mail: bizdjc@nus.edu.sg. The author thanks Yu-Hung Chien, Shuping Li, Kaican Kang and Qiqi Zou for their able research assistance.

# 1 Introduction

Selecting a suitable subset of regressors in the context of linear regressions or other similarly structured variable selection problems plays a critical role for a wide range of practical issues in the era of big data. Its importance is completely obvious so that we will skip the vast literature along many lines of research. Variable selection based on the zero-norm penalty (i.e., number of selected regressors) is conceptually more appealing than other criteria such as the  $l_1$ -norm penalty because it directly addresses the variable selection problem. Practically speaking, it also works better because regression coefficients are not distorted by the penalty term (i.e., shrinkage toward zero when being selected). Also interesting to note is the fact that the regression model fit, measured in  $R^2$ , is invariant to linearly transforming a group of regressors, but the corresponding  $l_p$  ( $0 < p < 2$ ) penalty term is not.<sup>1</sup> Therefore, multicollinearity which naturally occurs in data will interfere with regressor selection based on an  $l_p$  ( $0 < p < 2$ ) penalty, but not with the zero-norm regressor selection. In this paper, we propose a practical sequential Monte Carlo solution to the zero-norm regressor selection problem, and demonstrate through a simulation study the distortion caused by the  $l_1$ -norm based method such as Lasso.

The main idea of our zero-norm regressor selection method is to cast this variable selection problem as a pure combinatorial optimization and solve it with the recently emerged density-tempered sequential Monte Carlo (SMC) sampling technique by Del Moral, *et al* (2006), Duan and Fulop (2015), among others. Once a random combination of a fixed number of regressors is given, we rely on the analytical linear regression solution to attach to the combination a likelihood of occurrence, which in turn defines a probability distribution function, up to a missing normalizing constant, over all possible combinations of the same fixed number of variables. By sequential probability-tempering, one will arrive at a final SMC sample mimicking this target probability function whose maximum in turn provides the best combination at the fixed number of variables. Employing a set of, say, 1,000 SMC particles, our variable selection method is in essence a global solution approach.

Our SMC sampling technique originates from the Bayesian literature, but our approach differs from a long line of research papers using the Bayesian statistical techniques. Typical approaches rely on a hierarchical structure whether any given regression coefficient is modeled by a spike-and-slab (i.e., Bernoulli and Gaussian) mixture and its variants, for example, Mitchell and Beauchamp (1988), George and McCulloch (1993) and Polson and Sun (2018), or adaptive sampling for variables' inclusion via some parametric distribution for binary variables as in Schafer and Chopin (2013). In a true Bayesian sense, these algorithms typically

---

<sup>1</sup> $l_p$  for  $0 \leq p < 1$  is actually not a norm, and zero-norm so-named by David Donoho, a special case of  $p = 0$ , obviously lacks homogeneity. However, zero-norm has become a standard way of describing such a penalty form in the data analytics community.

require a strong prior and some further distribution assumptions to work, and essentially alters the original variable selection problem. In contrast, our SMC variable selection method does not make any additional distributional assumption and solves the combinatory problem in its original form. Practically important is our method’s ability to scale up for real big-data problems for which the number of potential regressors may run into thousands.

In a general spirit, our approach resembles that of Duan and Zhang (2016) where they generate a sample of non-Gaussian bridge paths by sequentially replacing random segments of the paths with suitable probabilities representing various intermediate target bridges, and eventually reach the final sample for the target non-Gaussian bridge model. We first show that selecting regressors subject to a zero-norm condition is equivalent to finding the maximum value of a distribution function defined over combinations of regressors with some fixed number of elements, say,  $k$ . Our target distribution function has no tractable analytical solution, but can be represented by a sample of  $k$ -dimensional points with each representing a  $k$ -combination of all potential regressors. Since different orders of regressors for a given combination yields the same regression solution, permutation has no particular meaning above and beyond defining a combination. However, a distribution function defined over permutations will be easier to work with, because permutations are easier to sample and the distribution function defined over permutations can be proportionally scaled up to obtain the target distribution function defined over combinations. Note that the SMC algorithm used in solving this maximization problem is by design free of any scaling constant. Hence, the SMC sample representing the distribution function over combinations is exactly the same as the one over permutations after collapsing permutations into combinations. On the methodological front, we devise a sensible but arbitrary initialization sampler to generate permutations and absorb its initialization distribution into the importance weight. Then, the progression of the algorithm is performed repeatedly through reweighting, resampling, and support-boosting steps to finally arrive at the SMC sample for the target distribution. For the Metropolis-Hastings move used in the support-boosting step, we also engage a proposal sampler defined over permutations for the same reason.

We conduct a comprehensive simulation study of selecting 9 or 18 variables out of 900 potential variables, and under which the number of potential combinations equals  $1.026 \times 10^{21}$  and  $1.986 \times 10^{37}$ , respectively. In the great majority of 500 simulation repetitions, the selected model has a higher in-sample  $R^2$  than that of the estimated true model, suggesting that the zero-norm SMC method can reliably select the best variable combination in a practical sense. The method’s performance improves when the observations are smaller and/or the residual errors are larger because sampling errors become larger and making room for the selected model to outperform the true model based on the in-sample  $R^2$ . Naturally, it also performs better when the number of selected variables is smaller, i.e., 9 versus 18 because the zero-norm SMC method is less likely to miss the best solution.

Another simulation study pertains to the relative performance of the zero-norm SMC method versus the Adaptive Lasso of Zou (2006). In this simulation study, the true model has 12 variables and the selection task is to choose the best combination among 900 potential variables with 200 observations while the number of variables is unknown to the analyst. Two-fold cross-validation is deployed to determine the number of variables in the final solution for both the zero-norm SMC method and the adaptive Lasso. The results clearly indicate that the adaptive Lasso is prone to pick too many variables with an average of 30 variables when the true model has 12 by design. Moreover, the range is quite wide with the solutions over 500 simulation repetitions ranging from 9 to 93. In contrast, the zero-norm SMC method yields an average of 8 variables and covers the range from 5 to 12 selected variables. Under-selection is actually expected because cross-validation is a conservative way of finding the “right” number of regressors by avoiding in-sample over-fitting in the case of 200 observations. The F-scores, a standard way of comparing models in machine learning, also suggest a far superior performance of the zero-norm SMC method.

The zero-SMC algorithm can in theory find the right solution with a probability of one when one increases the number of particle to infinity. Other than a theoretical interest, that would be no better than the brute force approach of exhausting all possible combinations. Similar to using the Central Limit Theorem to assess the Monte Carlo solutions in many contexts, our SMC solution’s  $R^2$  is the maximum order statistic of the final SMC sample and can thus be appraised with an estimated Weibull distribution, a max-stable distribution limit implied by the Fisher-Tippett-Gnedenko Extreme Value Theorem. We show by simulation that the predictive  $R^2$  distribution based on 2,000 SMC particles can well describe the likelihood and the extent of possible improvement upon the currently obtained maximum  $R^2$ .

At 2,000 SMC particles, a typical modern desktop computer can complete one selection task (for example, choosing 12 variables out of 900) in about one minute. Using cross-validation or the BIC criterion to determine the right number of selected variables, one can expect to complete the overall variable selection task under 30 minutes. In fact, one can push the number of particle to 10,000 or somewhat higher without running into a memory problem, and complete the selection task in a couple of hours.

## 2 Linear regression subject to a zero-norm penalty

Consider the classical linear regression model of  $p$  regressors with  $n$  observations:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \tag{1}$$

where  $\mathbf{y} = (y_1, \dots, y_n)'$ , and  $\mathbf{X}$  denotes the  $n$  observations of  $p$  regressors, i.e.,  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$  with  $\mathbf{x}_i = (x_{i1}, \dots, x_{in})'$ , of which the first vector may represent the intercept

term.  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$  is the  $p$ -dimensional regression coefficients, and  $\boldsymbol{\epsilon}$  is  $n$ -dimensional *i.i.d.* normally distributed errors with mean 0 and variance  $\sigma^2$ .

Sometimes, there are more potential regressors than data; that is,  $p$  is greater than  $n$ . Penalized regression is the only sensible way to estimate this regression. Even if there are enough data points, in-sample over-fitting is still a general modeling concern and penalized regression can be very useful in dealing with the over-fitting problem. Multicollinearity is another commonly encountered situation in regressions, meaning that some of them are highly correlated. Trimming away some regressors seems to be a wise thing to do, and penalized regression is an obvious way to go. Whatever is the reason, the general issue boils down to selecting a subset of good regressors that delivers a robust and reliable performance.

The penalized regression considered in this paper is the one subject to the zero-norm regularization.

$$\begin{aligned} \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{l_2}^2 \\ \text{s.t. } \|\boldsymbol{\beta}\|_{l_0} \leq p_s \leq p \end{aligned} \quad (2)$$

where  $\|\cdot\|_{l_2}$  and  $\|\cdot\|_{l_0}$  stands for the  $l_2$  and zero norms, respectively. Note that  $\|\boldsymbol{\beta}\|_{l_0}$  counts the number of non-zero entries in  $\boldsymbol{\beta}$ . Note that the above minimization problem is equivalent to  $\arg \min_{\boldsymbol{\beta}} \{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{l_2}^2 + \lambda\|\boldsymbol{\beta}\|_{l_0}\}$  where the solution is a step function of  $\lambda$  with the jumps corresponding to different values of  $k$ . This zero-norm penalized regression problem, albeit being natural for regressor selection, is known by Natarajan (1995) to be NP-hard. Thus, the Lasso technique of Tibshirani (1996) (i.e., replacing  $\lambda\|\boldsymbol{\beta}\|_{l_0}$  by  $\lambda\|\boldsymbol{\beta}\|_{l_1}$ ) has become extremely popular, because the minimization problem is convex and there is an efficient algorithm to find the unique global solution. However, the Lasso technique does not possess the oracle properties as defined by Fan and Li (2001). Alternatives to the Lasso and with the oracle properties are also popular, for example, the SCAD method of Fan (1997) and Fan and Li (2001) and adaptive Lasso of Zou (2006).

The ‘‘Irrepresentable Condition’’ of Zhao and Yu (2006) is an important factor in applying Lasso and its variants. The condition states that ‘‘Lasso selects the true model consistently if and (almost) only if the predictors that are not in the true model are ‘irrepresentable’ by predictors that are in the true model.’’ What it means in practice is that multicollinearity can create problems in applying Lasso or its variants. Many real problems come with many closely related variables. Individually, they offer meaningful explanatory powers, but including them all does not make sense nor actually improve the model’s performance after factoring in sampling errors. The practical issue is to find one or two of them work the best along with other variables of interest. Multicollinearity among these variables can be conceptually viewed as a natural result of linearly transforming a set of independent variables. Note that a regression model fit, measured in  $R^2$ , is invariant to such a linear transformation

of regressors, but the corresponding  $l_1$  penalty term is not. Such a natural set of competing variables may incur a high  $l_1$  penalty without measurably increasing the  $R^2$ , which then leads to their total elimination from the selected model. In short, multicollinearity interferes with regressor selection, and the Lasso or its variants may lead to the entire set of closely related variables not being selected instead of intuitively keeping one variable or two in the final model. This observation is not purely theoretical because we have encountered it often in practice.

In the next section, we devise a stochastic solution that is stable and runs efficiently. Theoretically, this stochastic solution can be made arbitrarily close to the true solution through the argument of the maximum order statistic. The solution technique also provides a natural way to trim or add regressors without having to restart from the scratch. By nature, our stochastic solution is a global optimization technique which is not plagued by short-sightedness or non-recoverability of the heuristic greedy stepwise algorithm often employed in regressor selection.

### 3 The distribution-tempered SMC solution

Our approach relies on reformulating the zero-norm constrained variable selection problem in (2) into an equivalent maximization problem over a distribution function defined on a  $p_s$ -dimensional random vector,  $\mathbf{U} = (U_1, U_2, \dots, U_{p_s})$ , that takes values from the set of the regressor sequence numbers  $\mathbf{P} = \{1, 2, \dots, p\}$  without replacement and its distribution is unknown. Naturally,  $p_s < n$  is understood. If otherwise, the regression model based on any combination of  $p_s$  regressors would, generally speaking, produce a perfect fit. We will leave the task of identifying its distribution function for a later discussion. For now, the focus is on setting up this equivalent reformulation. Without loss of generality, first note that the inequality constraint in (2) can be turned into an equality, because the solution will obviously be at the boundary. Thus, the minimization problem in (2) is equivalent to:

$$\arg \max_{\{\mathbf{U} \in \mathbf{P}(p_s)\}} \exp \left\{ -\|\mathbf{y} - \mathbf{X}_{\mathbf{U}}\hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2 \right\} \quad (3)$$

where  $\mathbf{P}(p_s)$  denotes  $\{\mathbf{U} \in \mathbf{P}^{p_s} \& U_1 \neq U_2 \neq \dots \neq U_{p_s}\}$ ,  $\mathbf{P}^{p_s}$  stands for the  $p_s$ -Cartesian product of  $\mathbf{P}$ ,  $\mathbf{X}_{\mathbf{U}}$  denotes the sub-matrix of  $\mathbf{X}$  whose columns correspond to the regressor sequence numbers in  $\mathbf{U}$ , and  $\hat{\boldsymbol{\beta}}(\mathbf{U}) = (\mathbf{X}_{\mathbf{U}}^t \mathbf{X}_{\mathbf{U}})^{-1} \mathbf{X}_{\mathbf{U}}^t \mathbf{y}$  is the optimal regression  $\boldsymbol{\beta}$  when  $\mathbf{U}$  is known.

Since  $\exp \left\{ -\|\mathbf{y} - \mathbf{X}_{\mathbf{U}}\hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2 \right\}$  is positive, it can obviously be viewed as a discrete distribution function of  $\mathbf{U}$ , up to a norming constant, over  $\mathbf{P}(p_s)$ , i.e, the set of permutations. Since the order of regressors for a given combination yields the same regression solution,

the distribution defined over permutations is a proportionally scaled down version of the distribution function defined over combinations where the scaling factor is  $p_s!$ . Since permutations are easier to generate, we will target the distribution function over permutations where the SMC algorithm is by design free of any proportional constant. Our target discrete distribution function is

$$f(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X}) \propto \exp \left\{ -\|\mathbf{y} - \mathbf{X}_{\mathbf{U}}\hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2 \right\}. \quad (4)$$

The task of selecting regressors has now been converted into finding the maximum of the above distribution function.

The basic idea of finding the maximum of (4) is to generate a sample suitably representing this distribution function. Our distribution-tempered SMC approach to finding this maximum comes from Del Moral, *et al* (2006) and Duan and Fulop (2015). Our approach resembles more that of Duan and Zhang (2016), which devises a way to effectively generate a high-dimensional random object subject to some condition. Generating a sample of non-Gaussian bridge paths is the target in Duan and Zhang (2016), whereas in this paper, the task is to find the maximum value of  $f(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X})$  through sequentially generating  $\mathbf{U}$  by tempering distribution to arrive at a sample that properly represents  $f(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X})$ . Thus, its maximal value and corresponding maximizer, up to a Monte Carlo error, becomes readily available.

### 3.1 The algorithm

Our distribution-tempered SMC method can be divided into three key steps - (1) initialization, (2) reweighting and resampling, and (3) support boosting.

#### Initialization

Assign each regressor with an initial probability  $\bar{q}_i$  for  $i \in \mathbf{P}$ . An intuitive and quick way is to set  $\bar{q}_i = R_i^2 / \sum_{j=1}^p R_j^2$  where  $R_i^2$  is the regression  $R^2$  using a single Variable  $i$ . Let  $q_i^{(p_s)}(0) = \bar{q}_i$  for all  $i$ 's, and use  $q_i^{(p_s)}(0)$  to sample  $\mathbf{U}$  from  $\mathbf{P}$ , i.e., simulate  $p_s$  regressors without replacement from the pool of  $p$  regressors. Since we need to obtain exactly  $p_s$  regressors, sampling can be performed sequentially by choosing the first regressor out of the  $p$  potential regressors with the probability of  $q_i^{(p_s)}(0)$  for  $i \in \mathbf{P}$ . Then, move on to selecting the second one out of the remaining  $p - 1$  regressors. Assuming that the first one is Variable  $i$ , the probability for choosing Variable  $j$  as the second becomes  $\frac{q_j^{(p_s)}(0)}{1 - q_i^{(p_s)}(0)}$  for  $j \in \mathbf{P} \setminus \{i\}$ . The same logic applies to the third regressor and so on until reaching the last one in the permutation, i.e.,  $p_s$ , for which the probability is  $q_k^{(p_s)}(0)$  for sampling Variable  $k$  from the remainder. If  $q_i^{(p_s)}(0) = 0$  for some Variable  $i$ , it will never be sampled and can thus be

trimmed from the set of regressors from the start. Thus,  $\mathbf{P}$  should be understood as the set of regressors with  $q_i^{(p_s)}(0) > 0$ . Use  $I(\mathbf{U} \in \mathbf{P}(p_s); q_i^{(p_s)}(0), i \in \mathbf{P})$  to denote this permutation sampler based on the distribution  $q_i^{(p_s)}(0)$ .

Note that this sampler's probability distribution, being a product of the probabilities just mentioned, depends on the specific sequence of appearance. In other words, the sampled points represent different permutations, which are easier to sample and evaluate their probabilities. The regression solution, however, only depends on the combination, and hence different permutations yielding the same combination will share the same value of  $\exp\{-\|\mathbf{y} - \mathbf{X}_U \hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2\}$ . Thus, the maximum is not unique over  $\mathbf{U} \in \mathbf{P}(p_s)$ , but it does not matter to the solution.

### Reweighting and resampling

Define an intermediate target distribution function as

$$f_\gamma(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X}) \propto \left( \frac{\exp\{-\|\mathbf{y} - \mathbf{X}_U \hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2\}}{I(\mathbf{U} \in \mathbf{P}(p_s); q_i^{(p_s)}(0), i \in \mathbf{P})} \right)^\gamma I(\mathbf{U} \in \mathbf{P}(p_s); q_i^{(p_s)}(0), i \in \mathbf{P}) \quad (5)$$

Obviously,  $f_\gamma(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X})$  equals  $I(\mathbf{U} \in \mathbf{P}(p_s); q_i^{(p_s)}(0), i \in \mathbf{P})$  when  $\gamma = 0$ , and  $f(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X})$  when  $\gamma = 1$ .

The distribution-tempered SMC technique moves the sample along a self-adapted control bridge by advancing  $\gamma$  from 0 to 1. The self-adapted control rests with choosing  $\gamma$  so that the effective sample size (ESS) implied by the importance weight does not fall below a threshold  $\eta M$  where  $M$  is the intended size of the SMC sample and  $\eta$  equals, say, 1/2. Denote

$$\text{the incremental importance weight by } w_{\gamma, \gamma^{(j)}}(\mathbf{U}^{(i)}) = \left( \frac{\exp\{-\|\mathbf{y} - \mathbf{X}_{\mathbf{U}^{(i)}} \hat{\boldsymbol{\beta}}(\mathbf{U}^{(i)})\|_{l_2}^2\}}{I(\mathbf{U}^{(i)} \in \mathbf{P}(p_s); q_i^{(p_s)}(0), i \in \mathbf{P})} \right)^{\gamma - \gamma^{(j)}} \quad \text{and}$$

$$\text{ESS} = \frac{(\sum_{i=1}^M w_{\gamma, \gamma^{(j)}}(\mathbf{U}^{(i)}))^2}{\sum_{i=1}^M (w_{\gamma, \gamma^{(j)}}(\mathbf{U}^{(i)}))^2}.$$

Set  $\gamma^{(0)} = 0$ . Find  $\gamma^*$  such that the ESS is no less than  $\eta M$ . Note that this solution need not be exactly at  $\eta M$ , because it is just a control device to prevent the quality of sample from deteriorating too much. Use the incremental importance weight to resample in order to obtain an equally-weighted sample of  $\mathbf{U}$ . Then, set  $\gamma^{(1)} = \gamma^*$ .

### Support boosting

After resampling, the sample is likely to contain more duplicate copies of some  $\mathbf{U}$ 's to reflect their relatively high importance weights, which means that empirical support has



shrunk.<sup>2</sup> We need to boost the empirical support before advancing  $\gamma$  again. Support boosting can be accomplished by several Metropolis-Hastings (MH) moves until the cumulative realized acceptance rate has reached a target level, say, 100%, which is to ensure that the empirical support has been properly boosted.<sup>3</sup>

The SMC sample provides a natural basis for coming up with a good proposal for executing MH moves. Compute  $c_i^{(p_s)}(\gamma) = \sum_{j=1}^M \sum_{l=1}^{p_s} \chi_{\{U_l^{(j)}=i\}}$ , which is the total count of Variable  $i$  appearing in the sample of size  $M$ . Define a probability by  $q_i^{(p_s)}(\gamma) = c_i^{(p_s)}(\gamma) / \sum_{j=1}^p c_j^{(p_s)}(\gamma)$  for  $i = 1, 2, \dots, p$ , which reflects the relative importance of Variable  $i$  after the SMC algorithm has reached the stage indicated by  $\gamma$ . If  $p_s$  regressors attain a probability of  $1/p_s$ , the sampler based on these probabilities will always generate different permutations of the same  $p_s$  regressors. We use  $\mathcal{Q}_U(\gamma)$  to represent these probabilities inferred from the sample. Instead of proposing a new permutation hoping to improve the solution, a more efficient and realistic approach is to replace a randomly selected subset of the existing permutation when  $p_s$  is large. Importantly, this differs from the single best replacement algorithm of Soussen, *et al* (2015).

Denote by  $h(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma))$  the conditional distribution based on  $\mathcal{Q}_U(\gamma)$  for proposing regressors to replace  $\tilde{A}$ , a random subject of elements in permutation  $\mathbf{U}$ , where  $\mathbf{U}_{-\tilde{A}}$  stands for removing from  $\mathbf{U}$  those elements in  $\tilde{A}$ . Given  $\tilde{A}$ , this probability is the standard permutation result for proposing element  $i \in \mathbf{P} \setminus \mathbf{U}_{-\tilde{A}}$ . Coupling with the probability of sampling  $\tilde{A}$  gives rise to the overall proposal probability.

One may be tempted to propose a replacement subset in a way just like the initialization step. However, if for some random chance the current  $M$   $p_s$ -variable combinations had completely missed the variables in the optimal solution, the support boosting step would forever miss the optimal choice. A more reliable way of proposing new regressor permutations is therefore to mix  $h(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma))$  with  $I(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \bar{q}_i, i \in \mathbf{P})$ , the initialization sampler applying to  $i \in \mathbf{P} \setminus \mathbf{U}_{-\tilde{A}}$ . Specifically, our proposal sampler is based on  $h^{(\omega)}(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma)) = \omega h(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma)) + (1 - \omega)I(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \bar{q}_i, i \in \mathbf{P})$ , where  $\omega$  is set to 50% in our implementation.

---

<sup>2</sup>It is worth noting that duplicates can be expected under a discrete distribution even for an ideal sample. Duplicates merely reflect the number of elements in the theoretical support versus the sample size. Support-boosting is meant to remove duplicates due to resampling, but it is impossible to get rid of duplicates inherent to a discrete distribution.

<sup>3</sup>Since the underlying distribution is discrete, complete distinctiveness of particles cannot be expected. To ensure that proper support boosting has been completed, one can attach to each particle a uniform random number whenever a proposal is made. Resampling destroys particle distinctiveness, but accepting new proposal has in effect restored distinctiveness which is revealed in the distinctiveness of these attached uniform random numbers.

The MH acceptance probability for replacing  $\tilde{A}$ , a random subset of elements in permutation  $\mathbf{U}$ , is

$$\begin{aligned}
& \alpha_\gamma^{(j)} \{ \mathbf{U} \in \mathbf{P}(p_s) \Rightarrow \mathbf{U}^* \in \mathbf{P}(p_s) \} \\
&= \min \left\{ 1, \frac{f_\gamma(\mathbf{U}^* \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X})}{f_\gamma(\mathbf{U} \in \mathbf{P}(p_s); \mathbf{y}, \mathbf{X})} \frac{h^{(\omega)}(\mathbf{U} \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma))}{h^{(\omega)}(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma))} \right\} \\
&= \min \left\{ 1, \frac{\exp \left\{ -\gamma \|\mathbf{y} - \mathbf{X}_{\mathbf{U}^*} \hat{\boldsymbol{\beta}}(\mathbf{U}^*)\|_{l_2}^2 \right\}}{\exp \left\{ -\gamma \|\mathbf{y} - \mathbf{X}_U \hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2 \right\}} \left( \frac{I(\mathbf{U}^* \in \mathbf{P}(p_s); \bar{q}_i, i \in \mathbf{P})}{I(\mathbf{U} \in \mathbf{P}(p_s); \bar{q}_i, i \in \mathbf{P})} \right)^{1-\gamma} \right. \\
&\quad \left. \times \frac{h^{(\omega)}(\mathbf{U} \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma))}{h^{(\omega)}(\mathbf{U}^* \in \mathbf{P}(p_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_U(\gamma))} \right\} \tag{6}
\end{aligned}$$

The acceptance probability defines a Markov kernel, and the target intermediate distribution in (5) is, by the standard argument, the stationary solution to the Markov kernel. In one round of support boosting, the MH move may replace a random number of elements to increase combination diversity. An average acceptance rate can be computed for a particular round. The support boosting step is considered satisfactorily completed when the cumulative average acceptance rate over rounds has reached the target level, and for which we set it at 500%.

After the empirical support is properly boosted, the system is ready for advancement from  $\gamma^{(1)}$  to  $\gamma^{(2)}$ , and eventually to 1. This is conducted by repeating the reweighting, resampling and support boosting steps.

Finally, the chosen  $p_s$  regressors corresponds to the point of the final sample that yields the highest value of the distribution function. As stated earlier, our sample points are permutations and the regression solution only depends on combinations. Hence, there are likely multiple SMC sample points attaining the same highest value of the distribution function, but they all amount to the same choice of  $p_s$  variables.

### 3.2 $k$ -fold duplication to enlarge the SMC sample

When the number of chosen regressors,  $p_s \leq p$ , is large, the size of the SMC sample,  $M$ , may need to be large enough to properly represent the underlying distribution over the discrete and yet very large set,  $\mathbf{P}(p_s)$ . In fact, if  $M$  is made arbitrarily large, the SMC solution will be arbitrarily close to the true solution in probability, or may even be the true solution simply because the underlying distribution is discrete with a large but finite number of combinations. Efficiently increasing the SMC sample size will thus make the method practically powerful. Duan and Zhang' (2016)  $k$ -fold duplication technique serves this purpose well.

$k$ -fold duplication is to duplicate the sample of size  $M$  to  $kM$  by making additional  $(k - 1)$  identical copies of the SMC sample of size  $M$ . One then relies on support-boosting to reduce duplicates (i.e., enlarge the empirical support) to turn the sample into a truly representative sample of size  $kM$ . The key to  $k$ -fold duplication vis-a-vis straightforward SMC with a sample of size  $kM$  is to directly leverage the final SMC sample of size  $M$  at the stage of  $\gamma = 1$  so as to bypass the intermediate steps required for the distribution-tempering bridge for  $(k - 1)M$  sample points. In the results reported, we have applied one round of 2-fold duplication to increase from 1,000 to 2,000 SMC particles .

### 3.3 Adding or trimming regressors

The final SMC solution with  $p_s$  regressors, denoted by  $\hat{\mathbf{U}}_{p_s}$ , serves as a good basis for adding and trimming variable(s) under a new target number  $p'_s$  that is different from  $p_s$ . The modified initialization sampler adopted here is a mixture sampler assigning a probability of  $\omega'$  to an adding/trimming operation. If  $p'_s < p_s$ , one can remove  $p_s - p'_s$  variables with equal probabilities from  $\hat{\mathbf{U}}_{p_s}$ . When  $p'_s > p_s$ , we sample  $p'_s - p_s$  variables from the remainder set of the potential variables per usual to add to  $\hat{\mathbf{U}}_{p_s}$ , and the sampling is performed according to the original initialization probabilities described earlier. We denote the probability for this adding/trimming sampler by  $T(\mathbf{U} \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P})$ . The mixture sampler also assigns a probability of  $1 - \omega'$  to the original initialization sampling, which chooses  $p'_s$  variables directly from the whole set of potential variables based on the original initialization sampling distribution. This modified initialization sampler for selecting  $p'_s$  regressors can be represented by  $I^{(\omega')}(\mathbf{U} \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P}) = \omega' T(\mathbf{U} \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P}) + (1 - \omega') I(\mathbf{U} \in \mathbf{P}(p'_s); \bar{q}_i, i \in \mathbf{P})$  where  $\omega'$  is set to 0.1 in our implementation.

Accompanying this mixture initialization, we need to modify equation (5) to

$$f_\gamma(\mathbf{U} \in \mathbf{P}(p'_s); \mathbf{y}, \mathbf{X}) \propto \left( \frac{\exp \left\{ -\|\mathbf{y} - \mathbf{X}_U \hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2 \right\}}{I^{(\omega')}(\mathbf{U} \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P})} \right)^\gamma I^{(\omega')}(\mathbf{U} \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P}) \quad (7)$$

Support boosting will be conducted per usual, but the MH acceptance probability must be altered accordingly to reflect the change to a mixture initialization sampler; that is,

$$\begin{aligned} & \alpha_\gamma^{(j)} \{ \mathbf{U} \in \mathbf{P}(p'_s) \Rightarrow \mathbf{U}^* \in \mathbf{P}(p'_s) \} \\ = & \min \left\{ 1, \frac{\exp \left\{ -\gamma \|\mathbf{y} - \mathbf{X}_{U^*} \hat{\boldsymbol{\beta}}(\mathbf{U}^*)\|_{l_2}^2 \right\}}{\exp \left\{ -\gamma \|\mathbf{y} - \mathbf{X}_U \hat{\boldsymbol{\beta}}(\mathbf{U})\|_{l_2}^2 \right\}} \left( \frac{I^{(\omega')}(\mathbf{U}^* \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P})}{I^{(\omega')}(\mathbf{U} \in \mathbf{P}(p'_s); \hat{\mathbf{U}}_{p_s}, \bar{q}_i, i \in \mathbf{P})} \right)^{1-\gamma} \right. \\ & \left. \times \frac{h^{(\omega, \omega')}(\mathbf{U} \in \mathbf{P}(p'_s) \mid \mathbf{U}_{-\bar{A}}^* = \mathbf{U}_{-\bar{A}}; \mathcal{Q}_U(\gamma))}{h^{(\omega, \omega')}(\mathbf{U}^* \in \mathbf{P}(p'_s) \mid \mathbf{U}_{-\bar{A}}^* = \mathbf{U}_{-\bar{A}}; \mathcal{Q}_U(\gamma))} \right\} \quad (8) \end{aligned}$$

where  $h^{(\omega, \omega')}(\mathbf{U} \in \mathbf{P}(p'_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_{\mathbf{U}}(\gamma))$  is the same as  $h^{(\omega)}(\mathbf{U} \in \mathbf{P}(p'_s) \mid \mathbf{U}_{-\tilde{A}}^* = \mathbf{U}_{-\tilde{A}}; \mathcal{Q}_{\mathbf{U}}(\gamma))$  in equation (6) except that the initialization sampler is now a mixture sampler involving  $\omega'$ .

### 3.4 Determining the right number of regressors

Thus far, we have described a density-tempered SMC way of finding the optimal regressor combination and the associated parameter estimates under a fixed number of regressors. In practice, one still needs to determine the “right” number of regressors to avoid in-sample over-fitting due to deploying too many variables.

We can deploy  $m$ -fold cross-validation to determine the suitable number of regressors. Per usual, one randomly divides the data sample into  $m$  groups and consider  $p_s$  regressors. First, take out one group as the testing dataset and use the remaining  $m - 1$  groups together as the training sample. Then, apply the optimal  $p_s$  regressors obtained from the training sample along with their regression coefficients to the testing sample to compute its sum of squared residuals (SSE). Third, repeat it for each of the  $m$  groups taken out as the testing dataset and add up  $m$  SSEs to obtain the out-of-sample total sum of squared residuals. Finally, searching over different  $p_s$  to find its optimal value that corresponds to the lowest total SSE. Since  $m$ -fold cross-validation is accompanied by additional sampling errors arising from the random division of the total sample, we adopt a spline smoothing technique over  $p_s$  and its corresponding total SSE. For the simulation study reported later involving cross-validation, we use a cubic spline with three knots at 2, 9 and 20 for the smoothing operation.

Alternatively, one can rely on the BIC or AIC criterion to determine the number of regressors. Recall that  $\hat{\mathbf{U}}_{p_s}$  is the optimal regressor choice when  $p_s$  variables are allowed. In our case, the model with  $p_s$  regressors has  $p_s + 1$  parameters (i.e.,  $p_s$  regression coefficients plus the residual variance), and the BIC and AIC under normality and leaving out the irrelevant constant are

$$\text{BIC}(p_s) = n \ln \left( \frac{1}{n} \|\mathbf{y} - \mathbf{X}_{\hat{\mathbf{U}}_{p_s}} \hat{\boldsymbol{\beta}}(\hat{\mathbf{U}}_{p_s})\|_{l_2}^2 \right) + (p_s + 1) \ln n \quad (9)$$

$$\text{AIC}(p_s) = n \ln \left( \frac{1}{n} \|\mathbf{y} - \mathbf{X}_{\hat{\mathbf{U}}_{p_s}} \hat{\boldsymbol{\beta}}(\hat{\mathbf{U}}_{p_s})\|_{l_2}^2 \right) + 2(p_s + 1) \quad (10)$$

One needs to compute the  $\text{BIC}(p_s)$  and  $\text{AIC}(p_s)$  under different numbers of selected regressors, and the model with the lowest BIC (or AIC) value gives rise to the final choice. Our experience suggests that either BIC or AIC tends to select more regressors vis-a-vis cross-validation.

### 3.5 Assessing reliability of the algorithm

Like all Monte Carlo algorithms, our zero-norm SMC selection method is a stochastic scheme that will in principle select the best  $p_s$ -variable combination out of the potential variables that gives the maximum  $R^2$  if the SMC sample size,  $M$ , approaches infinity. Practically speaking, however, one would need a means to assess how close the final solution's  $R^2$  is to the true maximum under some finite  $M$ . In typical Monte Carlo analyses, the Central Limit Theorem serves as the basis for such an assessment. In our context, the Fisher-Tippett-Gnedenko Extreme Value Theorem provides the theoretical basis.

Let  $\bar{R}^2(p_s)$  be the regression  $R^2$  at our final selected set of  $p_s$  variables. It is important to note that  $\bar{R}^2(p_s)$  may be strictly larger than  $R^2(p_s; M) = \max \{R_i^2(p_s); i = 1, 2, \dots, M\}$ , the set of  $R^2$  corresponding to different  $p_s$ -variable combinations in the final SMC sample of size  $M$ , because  $\bar{R}^2(p_s)$  is the best result recorded throughout the whole sequential updating process. Denote by  $R_{max}^2(p_s)$  the theoretical true maximum  $R^2$  after exhausting all possible combination of  $p_s$  variables, and naturally,  $R^2(p_s; M) \leq \bar{R}^2(p_s) \leq R_{max}^2(p_s)$ . Our tasks are (1) to estimate  $R_{max}^2(p_s)$  so as to know the potential room for improvement by increasing the SMC sample size, and (2) to estimate the probability of making further improvement, i.e., the probability of exceeding  $\bar{R}^2(p_s)$ .

$R^2(p_s)$  can be viewed as a random variable resulting from random  $p_s$ -variable combinations. Although its distribution function, denoted by  $G(\cdot)$ , is an unknown discrete distribution. Obviously,  $R^2$  is bounded above by its theoretical maximum at  $R_{max}^2(p_s) \leq 1$ . In the neighborhood of  $R_{max}^2(p_s)$ , one can easily show that  $G(x)$  can be approximated by  $(R_{max}^2(p_s) - x)^\alpha L(\frac{1}{R_{max}^2(p_s) - x})$  where  $L(\cdot)$  is a slowly varying function converging to zero at infinity and  $\alpha$  some unknown positive constant, because  $G(x)$  is fundamentally a discrete distribution. Thus, the Fisher-Tippett-Gnedenko Extreme Value Theorem implies convergence to a Weibull distribution with the shape parameter  $\alpha$  and scale parameter 1; that is, for large  $M' < M$

$$Prob \left\{ \frac{R^2(p_s, M') - R_{max}^2(p_s)}{R_{max}^2(p_s) - G^{\leftarrow}(1 - 1/M')} \leq x \right\} \cong \exp(-|x|^\alpha) \quad \text{for } x \leq 0 \quad (11)$$

where  $G^{\leftarrow}(x) \equiv \inf\{y : G(y) > x\}$ , i.e., the left continuous inverse. Consequently,  $R^2(p_s; M')$  has an approximate distribution function:

$$F_{R^2(p_s, M')}(z) = \exp \left[ - \left( \frac{R_{max}^2(p_s) - z}{\eta} \right)^\alpha \right] \quad \text{for } z \leq R_{max}^2(p_s) \quad (12)$$

where  $\eta = R_{max}^2(p_s) - G^{\leftarrow}(1 - 1/M')$ .

The Weibull distribution in (12) can be treated as a two- or three-parameter distribution function. Taking  $R_{max}^2(p_s)$ ,  $\alpha$  and  $\eta$  as unknown, the system has three parameters. If instead

we view  $G^{\leftarrow}(1 - 1/M')$  as known by its empirical distribution derived from the final SMC sample, then the system only has two unknown parameters, i.e.,  $R_{max}^2(p_s)$  and  $\alpha$ . We will estimate the unknown parameters using  $\{R_i^2(p_s); i = 1, 2, \dots, M\}$  by randomly partitioning the SMC sample into  $k$  subsamples of size  $M'$ ; for example, the SMC sample of size 2,000 are partitioned into 20 subsamples of 100 each. We then use these  $k$  subsample maximum  $R^2$ , i.e.,  $R^2(p_s; M')$  to find optimal values for  $R_{max}^2(p_s)$  ( $\leq 1$  and  $\geq \bar{R}^2(p_s)$ ),  $\alpha > 0$  and  $\eta > 0$  (if  $\eta$  is treated as a free parameter) by matching the  $k$ -point empirical distribution to the extreme value distribution by minimizing the  $l_2$  distance.<sup>4</sup>

Figure 1 displays a predicted Weibull distribution for  $R^2$  taken from a simulation study to be described in the next section. It is estimated to a 10-point empirical distribution deriving from 20 subsample maximum  $R^2$ 's where the data sample is generated with a model of 12 variables at a theoretical  $R^2$  of 80%, but the selection task is to choose 20 among 900 potential variables. In this figure,  $R_{max}^2 = 0.940329$  provides an estimate for how high the  $R^2$  can be potentially increased to from its current value of 0.937434, by increasing the number of SMC particles. Note that the grand maximum  $R^2$  (i.e., 0.939325) displayed on the graph is a proxy value for the maximum  $R^2$  under 20 variables obtained by running the SMC method numerous times with a random starting value. This example suggests that the estimated maximum  $R^2$  only slightly overstates the possibility.

## 4 A simulation study

To ascertain the reliability of the zero-norm SMC selection method, we conduct a simulation study and report the results in Table 1. This simulation study intends to address the question as to whether the best combination of explanatory variables can be found. Since the combinatory possibilities cannot be practically exhausted, there is no direct answer to the question, and we need an indirect and sensible way to assess the method's reliability. In this simulation study, we rely on the knowledge that the true model is always a feasible combination, and thus the estimated true model serves as a natural benchmark to assess the quality of our proposed zero-norm SMC method. Due to sampling errors, the in-sample  $R^2$  of the best combination must be greater than or equal to the benchmark value when the number of explanatory variables is kept the same as that of the true model. If the number of observations increases, it should become increasingly difficult for the best solution to yield an  $R^2$  strictly higher than the benchmark value.

---

<sup>4</sup>Note that some of these  $k$  subsample maxima may share common values. In which case, the  $l_2$  distance will be computed over fewer than  $k$  points. If the number of unique  $R^2$  values equals one, this degenerate case cannot be estimated and the true maximum  $R^2$  is considered already attained by  $\bar{R}^2(p_s)$ . When the number of unique  $R^2$  is two, we use the two-parameter distribution function whereas for cases with three or more unique  $R^2$  values, we always deploy the three-parameter distribution in estimation.

The simulation study reported in Table 1 always involves 900 potential explanatory variables, and we set out to select 9 or 18 variables, depending on which number were used to generate the data. The 900 variables are equally divided into three groups of 300 each. Within each group, the variables have the same correlation. The first group has zero correlation, the second is 0.4, and third is 0.8. Across groups, variables are independent. All 900 variables are normally distributed with mean 0 and variance 1. Among the 300 variables in each group, the regression coefficients are set to 0.1 for 100 variables, 0.5 for the second batch of 100 variables, and finally 1 for the remainder. We consider two simulation setups to ascertain the impact of residual errors in variable selection. We factor in the magnitude of residual errors by considering two levels of theoretical  $R^2$  at 80% and 40%. In addition, we examine the impact of the sample size and study the results under the sample size of 200 and 1,000, respectively. All studies are conducted with 500 repetitions to tally the various rates of occurrence.

We will use the top-left panel in Table 1 to explain the simulation results. This simulation study deploys 9 out of 900 variables to generate the data set as described above. The theoretical  $R^2$  is at 80% and each sample has 200 observations. The results show that the zero-norm SMC method is able to beat the estimated true model 100% out of 500 simulation repetitions. The results imply that the zero-norm SMC method has mostly likely found the best in-sample solution. The hit ratios (i.e, success in identifying the variables in the true model) reported in this panel for nine sub-categories suggest that the zero-norm SMC method yields a higher hit ratio for the lower correlation group and the variables with a larger regression coefficient. Low hit ratios for the high correlation group and the variables with a low regression coefficient are expected, because those variables are not supposed to be easily discernable when the sample size is 200. Moving on to other cases reported in different panels of Table 1 by varying the number of variables to be selected, the level of  $R^2$ , and sample size, a consistent pattern emerges to reflect the fact that sampling errors provide room for the zero-norm SMC method to outperform the estimated true model in-sample.

The comparison study between the zero-norm SMC method and the Adaptive Lasso of Zou (2006) is based on the true model having 12 variables among 900 potential variables. Each of 500 simulation repetitions generates 200 observations with the 900 variables again equally divided into three groups with different levels of within-group correlation of 0.1, 0.4 and 0.8. Each group is assigned four variables with regression coefficients of 0.1, 0.4, 0.7 and 1, respectively. The theoretical  $R^2$  is always fixed at 80%. Two-fold cross-validation is used to determine the “right” number of variables for both methods because in practice this number is unknown to the analyst. The results in Table 2 clearly indicate that the adaptive Lasso is prone to pick too many variables with an average of 30 variables when the target number is 12 by design. Moreover, the range is quite wide with the solutions over 500 simulation repetitions ranging from 9 to 93. In contrast, the zero-norm SMC method yields an average of

8 variables and covers the range from 5 to 12 selected variables. Under-selection is actually expected, because cross-validation is a conservative way of finding the “right” number of regressors by avoiding in-sample over-fitting in the case of 200 observations.

Three ratios reported in Table 2 need some explanations. First, the precision measures the number of selected variables among the 12 true variables in relation to the total number of variables being selected. A higher value implies a sort of higher accuracy. However, the precision can be misleading when a method under-selects, and this can be easily understood with an example. Suppose that a method only selects one variable which happens to be among the 12 true variables. The precision will be 100%, but the method has missed all 11 other true variables. The recall computes the number of selected variables among the 12 true variables divided by 12, the number of variables in the true model in this simulation study. This measure can also be quite misleading if a method over-selects. In the extreme case of selecting all variables, the recall will be by definition 100%. The F-score strikes a balance of the two, and is the harmonic average of precision and recall, i.e.,  $\left(\frac{\text{precision}^{-1} + \text{recall}^{-1}}{2}\right)^{-1}$ . Due to excessive over-selection by the adaptive Lasso, the recall is high and the precision is low, whereas the zero-norm SMC method goes in an opposite direction. All in all, the F-score of 0.7 for the zero-norm SMC method implies its superior performance over the adaptive Lasso’s 0.43.

The hit ratios for all 12 sub-categories in Table 1 are greater for the adaptive Lasso, reflecting the fact that it has greatly over-selected variables. The zero-norm SMC method tends to under-select when cross-validation is used as the criterion to determine the number of variables. Again, this is not inherent to the zero-norm SMC method, rather it reflects the conservative nature of cross-validation in avoiding in-sample over-fitting. For a variable with a small regression coefficient, sampling errors tend to bury its identity as a variable in the true model.

Our next simulation study is to determine how reliable the prediction based on the extreme value theorem described in Section 3.5 in predicting the true maximum  $R^2$  and generating a probability for further improvement upon the currently obtained maximum  $R^2$  in one sample. For this simulation study, we use one fixed simulated sample of 200 data points based on the 12-variable setup with 900 potential regressors in the simulation study reported in Table 2. We randomize the SMC variable selection algorithm, meaning that every variable selection run will begin with a random seed and repeat the SMC selection 500 times on the same data sample. Since the data set is fixed, the true maximum  $R^2$  is fixed. However, the value predicted by the extreme value theory will vary over the 500 runs. Although the true model is based 12 variables, analysts are not supposed to know this fact. So, we examine the performance under different numbers of selected variables ranging from 1 to 20, and in each case, we study how well the extreme value theory performs.



Each SMC run produces a target sample size of 2,000, which is always partitioned into 20 subsamples of 100 points each. We compute 20 subsample maxima and use them to estimate the Weibull distribution as described in Section 3.5. Since we do not know the true value of  $R_{max}^2(p_s)$ , it is reasonable to use a grand maximum over the 600  $\bar{R}^2(p_s)$ 's, which is obtained by having 100 simulation runs in addition to the original 500 runs. This grand maximum  $R^2$  is used to proxy the unknown true maximum  $R^2$  under each  $p_s$ . This design gives rise to a sample of 500 comparisons with each estimated  $R_{max}^2(p_s)$  against the grand maximum  $R^2$ . Each exceedance instance (i.e., the grand maximum  $R^2$  is strictly greater than an individual  $\bar{R}^2(p_s)$ ) over 500 simulation runs can be checked against whether the extreme value theory predicted exceedance probability is strictly positive, suggesting a possibility of improvement. The agreement rate can be tallied over the 500 simulation runs to serve as one way of gauging the performance of the extreme value theory in this application. Table 3 suggests a very high agreement rate with the worst case still having 88.2%.

The exceedance probability may suggest a possibility of  $R^2$  improvement, but the improvement magnitude may still be quite small if one substantially increases the SMC sample size. Table 4 shows the results by checking the difference between the grand  $R_{max}^2(p_s)$  described earlier and  $\bar{R}^2(p_s)$  of the SMC selected model. The various statistics on the differences over 500 random SMC runs on a same data sample of size 200 are tallied and reported in this table. Evidently, the improvement potential is relatively larger for a larger  $p_s$ . However, even when  $p_s = 20$ , the  $R^2$  improvement is on average less than 0.4% and the maximum improvement over the 500 runs is only slightly higher than 0.5%. In short, our proposed SMC method performs really well.

Next we check how close the predicted maximum  $R^2$  is to the true maximum  $R^2$  (proxied by the grand maximum  $R^2$ ) with the results reported in Table 5. The predicted value can over- or under-estimate the true value, and the magnitude of prediction error tends to increase with the number of selected variables, i.e.,  $p_s$ . These prediction errors are obviously small though with the maximum gap being just over 0.5%, suggesting the extreme value theory performs well in this application.

## References

- [1] Del Moral, P., A. Doucet, and A. Jasra, 2006, Sequential monte carlo samplers, *Journal of the Royal Statistical Society: Series B*, 68(3), 411-436.
- [2] Duan, J.-C., and A. Fulop, 2015, Density-Tempered Marginalized Sequential Monte Carlo Samplers, *Journal of Business and Economic Statistics*, 33(2), 192-202.
- [3] Duan, J.C. and C. Zhang, 2016, Non-Gaussian Bridge Sampling with an Application, National University of Singapore working paper.

- [4] Fan, J., 1997, Comments on “Wavelets in Statistics: a Review” by A. Antoniadis, *Journal of the Italian Statistical Society* 6(20), 131-138.
- [5] Fan, J. and R. Li, 2001, Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties, *Journal of the American Statistical Association* 96, 1348-1360.
- [6] George, E. and R. McCulloch, 1993, Variable Selection via Gibbs Sampling, *Journal of the American Statistical Association*, 88, 881-889.
- [7] Mitchell, T. and J. Beauchamp, 1988, Bayesian Variable Selection in Linear Regression (with discussion), *Journal of the American Statistical Association* 83, 1023-1036.
- [8] Natarajan, B.K., 1995, Sparse Approximate Solutions to Linear Systems, *SIAM Journal on Computing*, 24(2), 227-234.
- [9] Polson, N. and L. Sun, 2018, Bayesian  $l_0$ -regularized Least Squares, *Applied Stochastic Models in Business and Industry*, 1-15.
- [10] Schafer, C and N. Chopin, 2013, Sequential Monte Carlo on Large Binary Sampling Spaces, *Statistics and Computing* 23, 163-184.
- [11] Soussen C., J. Idier, J. Duan, and D. Brie, 2015, Homotopy Based Algorithms for  $l_0$ -regularized Least-squares, *IEEE Trans Signal Process*, 63, 3301-3316.
- [12] Tibshirani, R., 1996, Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society, Ser. B*, 58(1), 267-288.
- [13] Zhao, P. and B. Yu, 2016, On Model Selection Consistency of Lasso, *Journal of Machine Learning Research* 7, 2541-2563.
- [14] Zou, H., 2006, The Adaptive Lasso and Its Oracle Properties, *Journal of the American Statistical Association*, 101(476), 1418-1429.

Figure 1: An example of the extreme value theory predicted  $R^2$  distribution function

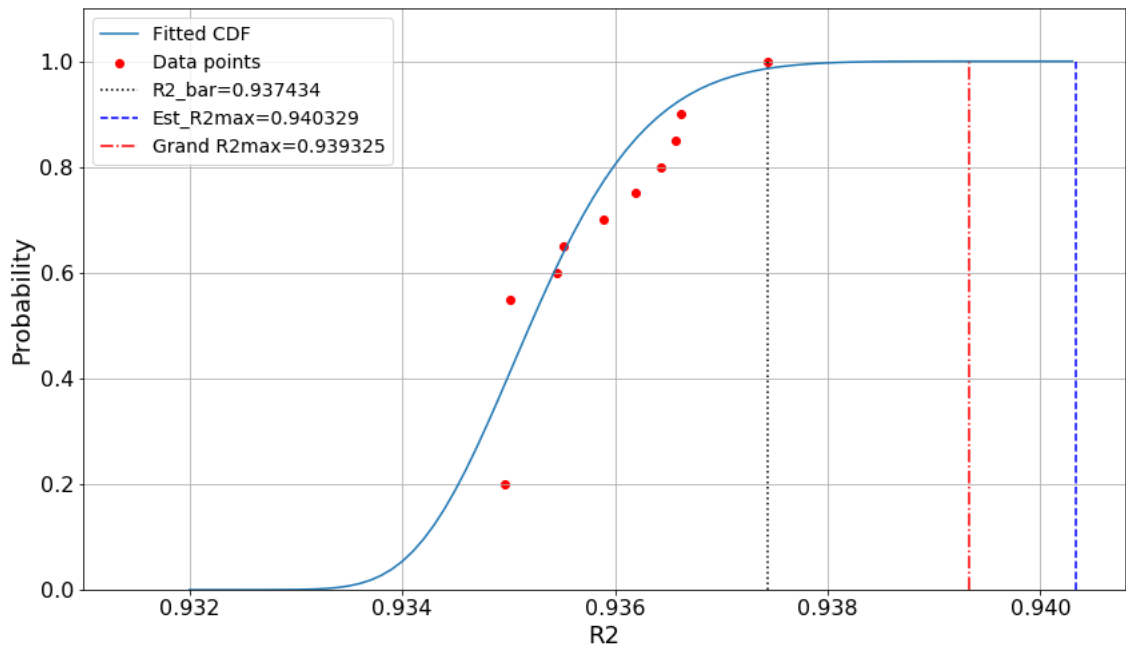


Table 1: A simulation study of zero-norm SMC variable selection out of 900 potential variables: the selected model vs. the estimated true model over 500 simulated samples. The 900 variables, which are normal random variables of mean 0 and variance 1, are divided into three independent groups of 300 each and at three levels of within-group correlation – 0, 0.4, and 0.8. The magnitude of regression coefficient varies from low to high and the theoretical  $R^2$  of the simulated model is at 40% or 80%.

<b># of obs=200, # of var=9, and <math>R^2=80\%</math></b>				<b># of obs=200, # of var=9, and <math>R^2=40\%</math></b>			
$R^2$ of the selected model is strictly greater: 100%				$R^2$ of the selected model is strictly greater: 100%			
$R^2$ of the selected model is greater or equal: 100%				$R^2$ of the selected model is greater or equal: 100%			
<b>Hit Ratio</b>		Within-group correlation		<b>Hit Ratio</b>		Within-group correlation	
Coefficient	0	0.4	0.8	Coefficient	0	0.4	0.8
0.1	0.048	0.024	0.018	0.1	0.008	0.008	0.014
0.5	0.998	0.988	0.592	0.5	0.454	0.250	0.076
1.0	1.0	1.0	0.992	1.0	0.992	0.934	0.394
<b># of obs=200, # of var=18, and <math>R^2=80\%</math></b>				<b># of obs=200, # of var=18, and <math>R^2=40\%</math></b>			
$R^2$ of the selected model is strictly greater: 99.8%				$R^2$ of the selected model is strictly greater: 100%			
$R^2$ of the selected model is greater or equal: 99.8%				$R^2$ of the selected model is greater or equal: 100%			
<b>Hit Ratio</b>		Within-group correlation		<b>Hit Ratio</b>		Within-group correlation	
Coefficient	0	0.4	0.8	Coefficient	0	0.4	0.8
0.1	0.01	0.022	0.018	0.1	0.004	0.012	0.026
0.2	0.1	0.116	0.056	0.2	0.036	0.036	0.030
0.5	0.832	0.79	0.276	0.5	0.180	0.144	0.058
0.7	0.984	0.986	0.576	0.7	0.364	0.248	0.122
0.85	0.992	0.998	0.774	0.85	0.566	0.436	0.160
1.0	1.0	1.0	0.89	1.0	0.748	0.544	0.256
<b># of obs=1000, # of var=9, and <math>R^2=80\%</math></b>				<b># of obs=1000, # of var=18, and <math>R^2=80\%</math></b>			
$R^2$ of the selected model is strictly greater: 95.8%				$R^2$ of the selected model is strictly greater: 94.6%			
$R^2$ of the selected model is greater or equal: 95.8%				$R^2$ of the selected model is greater or equal: 94.6%			
<b>Hit Ratio</b>		Within-group correlation		<b>Hit Ratio</b>		Within-group correlation	
Coefficient	0	0.4	0.8	Coefficient	0	0.4	0.8
0.1	0.318	0.344	0.074	0.1	0.106	0.166	0.032
0.5	1.0	1.0	1.0	0.2	0.606	0.784	0.244
1.0	1.0	1.0	1.0	0.5	1.0	1.0	0.992
				0.7	1.0	1.0	1.0
				0.85	1.0	1.0	1.0
				1.0	0.998	1.0	1.0

Table 2: A comparison study of zero-norm SMC variable selection vs. adaptive Lasso in selecting out of 900 potential variables over 500 simulated samples. The 900 variables, which are normal random variables of mean 0 and variance 1, are divided into three independent groups of 300 each and at three levels of within-group correlation – 0, 0.4, and 0.8. The true number of variables is set at 12 with four assigned to each group, and their coefficients are 0.1, 0.4, 0.7 and 1. The number of observations is fixed at 200 and the theoretical  $R^2$  of the simulated model is set to 80%. Two-fold cross-validation is deployed to determine the number of selected variables for both methods.

<b>Performance</b>	<b>Zero-norm SMC</b>				<b>Adaptive Lasso</b>			
# of selected variables (min, mean, max)	(5, 8, 12)				(9, 30, 93)			
F-Score	0.70				0.43			
Precision	0.88				0.32			
Recall	0.58				0.70			
<b>Hit Ratio</b>	Within-group correlation				Within-group correlation			
	Coef	0	0.4	0.8	Coef	0	0.4	0.8
	0.004	0	0.006	0.008	0.1	0.164	0.076	0.050
	0.4	0.545	0.541	0.161	0.4	0.954	0.880	0.446
	0.7	0.983	0.992	0.700	0.7	1.0	1.0	0.878
1.0	1.0	1.0	0.987	1.0	1.0	1.0	1.0	

Table 3: The extreme value theory predicted probability for  $R^2$  improvement under different  $p_s$  (the number of selected regressors) over 500 random SMC runs on a same data sample of size 200. The agreement rate is the rate of occurrence over 500 runs where the predicted probability is strictly positive and the grand  $R_{max}^2(p_s)$  (the proxy for the true maximum  $R^2$ ) is strictly greater than  $\bar{R}^2(p_s)$  (the  $R^2$  of the SMC selected model). The data sample is simulated with a generating model of 12 variables at a theoretical  $R^2$  of 80%. Four variables are assigned to each of the three correlation groups (i.e., 0, 0.4 and 0.8) with their coefficients equal to 0.1, 0.4, 0.7 and 1, respectively.

$p_s$	Average Exceedance Probability	Agreement Rate
1 to 7	0	1
8	0.000020	0.996
9	0.000001	0.998
10	0.000000	1
11	0.005427	0.882
12	0.003247	0.958
13	0.013839	0.890
14	0.020048	0.964
15	0.012050	1
16	0.029059	1
17	0.022858	1
18	0.029028	1
19	0.038637	0.998
20	0.036572	1

Table 4: Differences between the grand  $R_{max}^2(p_s)$  (the proxy for the true maximum  $R^2$ ) and  $\bar{R}^2(p_s)$  (the  $R^2$  of the SMC selected model) under different  $p_s$  (the number of selected regressors) over 500 random SMC runs on a same data sample of size 200. The data sample is simulated with a generating model that contains 12 variables with a theoretical  $R^2$  of 80%. Four variables are assigned to each of the three correlation groups (i.e., 0, 0.4 and 0.8) with their coefficients equal to 0.1, 0.4, 0.7 and 1, respectively.

$p_s$	Mean	Std	Min	25%	Median	75%	Max
1 to 6	0	0	0	0	0	0	0
7	0.001106	0.003382	0	0	0	0	0.011742
8	0.000132	0.001284	0	0	0	0	0.016442
9	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0
11	0.000054	0.000227	0	0	0	0	0.001497
12	0.000184	0.000334	0	0	0	0	0.001088
13	0.000178	0.000338	0	0	0	0.000180	0.001662
14	0.000697	0.000676	0	0	0.000342	0.001406	0.002961
15	0.001324	0.000624	0	0.000938	0.001075	0.001202	0.003179
16	0.002021	0.000494	0	0.001804	0.001869	0.002108	0.004019
17	0.002595	0.000492	0	0.002367	0.002724	0.002738	0.004343
18	0.003210	0.000585	0	0.002801	0.003412	0.003539	0.004860
19	0.003563	0.000704	0	0.003273	0.003740	0.003942	0.005527
20	0.003570	0.000763	0	0.003329	0.003792	0.004063	0.005462

Table 5: Differences between the grand  $R_{max}^2(p_s)$  (the proxy for the true maximum  $R^2$ ) and the extreme value theory predicted  $R_{max}^2(p_s)$  under different  $p_s$  (the number of selected regressors) over 500 random SMC runs on a same data sample of size 200. The data sample is simulated with a generating model that contains 12 variables with a theoretical  $R^2$  of 80%. Four variables are assigned to each of the three correlation groups (i.e., 0, 0.4 and 0.8) with their coefficients equal to 0.1, 0.4, 0.7 and 1, respectively.

$p_s$	Mean	Std	Min	25%	Median	75%	Max
1 to 6	0	0	0	0	0	0	0
7	0.001051	0.003306	0	0	0	0	0.011742
8	0.000142	0.001392	-0.002681	0	0	0	0.016442
9	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0
11	-0.000290	0.002459	-0.031069	0	0	0	0.001497
12	0.000023	0.001243	-0.017863	0	0	0	0.001088
13	-0.000373	0.002871	-0.024951	0	0	0	0.001549
14	-0.000212	0.004165	-0.045724	0	0.000342	0.001405	0.002502
15	0.000911	0.002411	-0.022633	0.000938	0.000938	0.001188	0.003179
16	0.001218	0.003104	-0.022319	0.001704	0.001804	0.002108	0.004010
17	0.001872	0.003148	-0.025740	0.002367	0.002524	0.002737	0.004343
18	0.002116	0.004321	-0.036267	0.002644	0.003412	0.003519	0.004768
19	0.002571	0.003387	-0.018868	0.003141	0.003709	0.003942	0.005084
20	0.002223	0.005220	-0.047293	0.003115	0.003609	0.004011	0.004726