

1

2 **Supplementary Information for**

3 **A Polynomial Algorithm for Best Subset Selection Problem**

4 **Junxian Zhu, Canhong Wen, Jin Zhu, Heping Zhang and Xueqin Wang**

5 **Heping Zhang and Xueqin Wang.**

6 **E-mail: heping.zhang@yale.edu; wangxq88@mail.sysu.edu.cn**

7 **This PDF file includes:**

8 Supplementary text

9 Figs. S1 to S2

10 Tables S1 to S3

11 SI References

12 Supporting Information Text

13 Computation

14 We build an R package called ABESS that implements our algorithms. Two major strategies to accelerate the computation are
15 noteworthy.

16 **Warm start** : A good initial active set can shorten the computation time. For ABESS, we can use the previous solution as an
17 initial set for the next step.

18 **Max splicing size** : Splicing takes the most computation in ABESS. It requires the order of k_{\max} operations in each
19 iteration. We set the max splicing size $k_{\max} = s$ in Algorithm 1. Based on our experience, we suggest to use a splicing size in
20 $k_{\max}, k_{\max}/2, \dots, 1$.

21 Additional simulation studies

22 **High-dimensional case.** Here, we intend to compare ABESS with more variable selection methods under the same setting in
23 our body content. The variable selection methods take into consideration include state-of-the-art variable selection algorithms
24 (i.e., LASSO, SCAD, and MCP) and recent proposed best subset selection solvers (i.e., SDAR, coordinate-wise minimizer,
25 and partial swap inescapable minimizer). Analogous to LASSO/SCAD/MCP, coordinate-wise minimizer (CWM) solves ℓ_0
26 regularization ordinary least squares estimator via the coordinate descent algorithm (1). In regard to partial swap inescapable
27 minimizer (PSIM) and SDAR, their solve best subset selection by local combinatorial search algorithm and primal-dual active
28 set algorithm (1, 2). We use their R implementations, *ncvreg*, *glmnet*, *L0Learn* and *BeSS* (3–6), in our simulation studies. As
29 for tuning parameters selection, we employ CV, SIC, EBIC and BIC to pick out the optimal one from l pre-specified tuning
30 parameters, which are default values given by their implementations. All methods to be compared are summarized in Table S1.
31 The synthetic datasets are generated following the same procedure in body content. The simulation result are exhibited in
32 Tables S2 and S3.

33 We first investigate the results in the uncorrelated-covariate regime (See Tables S2). In this regime, LASSO or the best
34 subset selection solvers (BSSS) are the best true positive detectors, and MCP or SCAD or information criterion (IC) based
35 BSSS have the least the false positive discovery. The notable advantage of BSSS is the balance of the true and false discoveries,
36 especially in $p \gg n$ case. As regard to parameter estimation, ABESS/SDAR/CWM/PSIM gives the most accurate estimations
37 among all of the methods, and specifically, the SDAR-EBIC have the most precise estimation when $p \leq 1500$ and ABESS-SIC
38 (or SDAR-SIC) takes the first place when p increases to 2500. Moreover, it can be seen that the the model size of the IC based
39 BSSS is closed to the ground truth. In terms of computational time, the most efficient method is IC based LASSO estimator,
40 and the second place is IC based SCAD/MCP, follow by ABESS, which is the fastest algorithms among all BSSS. Since CV
41 and local combinatorial search are two time-consuming techniques, PSIM and CV based methods generally require much more
42 time to conduct algorithms. It is also worthy to note that, in the uncorrelated case, ABESS, SDAR-SIC, CWM-SIC, PSIM-SIC
43 have a quite similar performance because global optimal, coordinate-wise optimal, (i.e., SDAR-SIC and CWM-SIC) and partial
44 swap inescapable optimal (PSIM) are likely to be equivalent in this scenario. Next, we study the simulation result under the
45 high-correlation regime (See Tables S3). By contrast with the uncorrelated case, both TPR and TNR of almost all methods
46 decrease in the context of high correlation, but LASSO/BSSS still have the highest TPR and MCP/SCAD/BSSS still have
47 the highest TNR. Notably, among SIC base BSSS, ABESS has a preferable variable selection in this scenario. As for the
48 parameter estimation, ABESS or SDAR-EBIC takes a dominant role due to the virtue of both algorithm and information
49 criterion. Besides, the best subset selection solvers equipped with SIC or EBIC are capable of detecting a proper model size.
50 Computationally, ABESS is the fastest best subset selection solvers and LASSO is the fastest variable selection solvers.

51 **Comparison with SDAR and IHT.** In this part, we study the convergence problems of SDAR and IHT algorithms.
52 Consider the classical linear model setting in our main body, the rows of $\mathbf{X}_{n \times p}$ are *i.i.d* sampled from the multivariate normal
53 distribution with covariance matrix Σ , where $\Sigma = (\sigma_{ij})_{p \times p}$ follows the three structures: uncorrelated (i.e., $\sigma_{ij} = \mathbf{I}(i = j)$),
54 decayed correlation (i.e., $\sigma_{ij} = 0.8^{|i-j|}$) and constant correlation (i.e., $\sigma_{ij} = 0.8^{\mathbf{I}(i \neq j)}$). For coefficient vector β , it has 10
55 non-zero elements, half of which are 1 and half -1. The n error terms are *i.i.d* drawn from the normal distribution $N(0, \sigma^2)$,
56 where $\sigma^2 = \beta^T \Sigma \beta$ such that the signal-noise ratio is fixed at 1.

57 We compare IHT, SDAR and ABESS for the ℓ_0 constraint problem with a fixed support size in terms of the number of
58 iterations. The maximum number of iterations is set to 100. We generate 100 synthetic datasets containing $n = 300$ observations
59 with $p = 1500$ covariates. The IHT algorithm is implemented in R and the step size in each iteration is determined by the
60 line search strategy. It stops when the ℓ_2 -norm of two consecutive estimations is smaller than 10^{-6} . We employ the SDAR
61 algorithm implemented in the *BeSS* package (2).

62 From Figure S1, the TPR and TNR of ABESS is significantly superior to the other algorithms in all of the settings, which
63 implies a better support recovery performance of ABESS, and in the meantime, ABESS also possesses the best estimation
64 performance in all situations. We note from the right-bottom panel of Figure S1 that SDAR may not converge in all of the cases
65 due to the potential periodic iterative issue, and even worse, it cannot converge with a high probability when the correlation
66 structure decays. For the IHT algorithm, its convergence rate declines as the global correlation increases. However, ABESS
67 gets rid of both problems and generally converges to the solution within 5 iterations.

68 **Faster sparsity level selection strategy.** Golden section (GS) search (11) is a strategy for finding an extremum of a
69 function inside a specified interval without regularity conditions such as continuity and derivative. It is a computationally
70 cheap strategy since it avoids some redundant computation steps. By applying GS on the search of the best sparsity level, we

Table S1. All methods considered in the simulation study .

Method	Problem	Algorithm	Selector	R-package	Reference
ABESS	Best subset selection	ABESS	SIC	splicing	this paper
SDAR-EBIC	Best subset selection	SDAR	EBIC	BeSS	(2)
SDAR-SIC	Best subset selection	SDAR	SIC	BeSS	(7)
CWM-CV	Best subset selection	CWM	CV	L0Learn	(1)
CWM-SIC	Best subset selection	CWM	SIC	L0Learn	this paper
PSIM-CV	Best subset selection	PSIM	CV	L0Learn	(1)
PSIM-SIC	Best subset selection	PSIM	SIC	L0Learn	this paper
MCP-CV	Noncave regularization	CWM	CV	ncvreg	(4)
MCP-SIC	Noncave regularization	CWM	SIC	ncvreg	(8)
MCP-BIC	Noncave regularization	CWM	BIC	ncvreg	(9)
MCP-EBIC	Noncave regularization	CWM	EBIC	ncvreg	(9)
SCAD-CV	Noncave regularization	CWM	CV	ncvreg	(4)
SCAD-SIC	Noncave regularization	CWM	SIC	ncvreg	(8)
SCAD-BIC	Noncave regularization	CWM	BIC	ncvreg	(9)
SCAD-EBIC	Noncave regularization	CWM	EBIC	ncvreg	(9)
LASSO-CV	ℓ_1 regularization	CWM	CV	glmnet	(3)
LASSO-SIC	ℓ_1 regularization	CWM	SIC	glmnet	this paper
LASSO-BIC	ℓ_1 regularization	CWM	BIC	glmnet	(10)
LASSO-EBIC	ℓ_1 regularization	CWM	EBIC	glmnet	(9)

Table S2. TPR, TNR, specificities, relative errors and sparse level errors of nineteen variable selection methods for uncorrelated correlation structure.

	Method	TPR	TNR	ReErr	SLE	Runtime
$p = 500$	ABESS	0.959 (0.064)	0.999 (0.001)	0.125 (0.112)	0.060 (0.908)	0.145 (0.029)
	SDAR-EBIC	0.959 (0.064)	0.999 (0.001)	0.110 (0.107)	-0.120 (0.795)	0.950 (0.166)
	SDAR-SIC	0.959 (0.064)	0.999 (0.001)	0.125 (0.112)	0.060 (0.908)	0.919 (0.167)
	CWM-CV	0.965 (0.059)	0.997 (0.003)	0.187 (0.173)	1.020 (1.531)	1.423 (0.209)
	CWM-SIC	0.959 (0.064)	0.999 (0.001)	0.125 (0.112)	0.060 (0.908)	0.176 (0.039)
	PSIM-CV	0.965 (0.059)	0.997 (0.002)	0.181 (0.165)	0.950 (1.359)	7.990 (1.186)
	PSIM-SIC	0.959 (0.064)	0.999 (0.001)	0.125 (0.112)	0.060 (0.908)	0.907 (0.174)
	MCP-CV	0.830 (0.137)	1.000 (0.000)	1.823 (1.598)	-1.700 (1.374)	0.530 (0.134)
	MCP-SIC	0.830 (0.137)	1.000 (0.000)	1.823 (1.598)	-1.700 (1.374)	0.074 (0.018)
	MCP-EBIC	0.830 (0.137)	1.000 (0.000)	1.823 (1.598)	-1.700 (1.374)	0.078 (0.022)
	MCP-BIC	0.830 (0.137)	1.000 (0.000)	1.823 (1.598)	-1.700 (1.374)	0.074 (0.018)
	SCAD-CV	0.831 (0.138)	1.000 (0.000)	3.108 (2.488)	-1.690 (1.383)	0.511 (0.120)
	SCAD-SIC	0.831 (0.138)	1.000 (0.000)	3.108 (2.488)	-1.690 (1.383)	0.076 (0.019)
	SCAD-EBIC	0.831 (0.138)	1.000 (0.000)	3.108 (2.488)	-1.690 (1.383)	0.077 (0.021)
	SCAD-BIC	0.831 (0.138)	1.000 (0.000)	3.108 (2.488)	-1.690 (1.383)	0.075 (0.019)
	LASSO-CV	0.968 (0.058)	0.973 (0.031)	0.534 (0.183)	12.710 (15.330)	0.446 (0.076)
LASSO-SIC	0.966 (0.057)	0.997 (0.003)	0.736 (0.476)	1.360 (1.685)	0.048 (0.023)	
LASSO-EBIC	0.966 (0.057)	0.996 (0.004)	0.712 (0.458)	1.660 (1.871)	0.047 (0.020)	
LASSO-BIC	0.969 (0.056)	0.991 (0.007)	0.601 (0.408)	4.150 (3.436)	0.048 (0.021)	
$p = 1500$	ABESS	0.972 (0.049)	1.000 (0.001)	0.154 (0.153)	0.260 (0.949)	0.342 (0.051)
	SDAR-EBIC	0.972 (0.049)	1.000 (0.001)	0.145 (0.151)	0.170 (0.900)	1.913 (0.337)
	SDAR-SIC	0.972 (0.049)	1.000 (0.001)	0.154 (0.153)	0.260 (0.949)	1.894 (0.332)
	CWM-CV	0.974 (0.046)	0.999 (0.001)	0.201 (0.157)	0.920 (1.116)	4.637 (0.611)
	CWM-SIC	0.972 (0.049)	1.000 (0.001)	0.155 (0.153)	0.290 (1.018)	0.536 (0.074)
	PSIM-CV	0.974 (0.046)	0.999 (0.001)	0.202 (0.159)	0.920 (1.116)	33.876 (4.131)
	PSIM-SIC	0.972 (0.049)	1.000 (0.001)	0.155 (0.153)	0.290 (1.018)	3.750 (0.636)
	MCP-CV	0.857 (0.128)	1.000 (0.000)	1.930 (1.539)	-1.430 (1.281)	1.440 (0.223)
	MCP-SIC	0.857 (0.128)	1.000 (0.000)	1.930 (1.539)	-1.430 (1.281)	0.224 (0.040)
	MCP-EBIC	0.857 (0.128)	1.000 (0.000)	1.930 (1.539)	-1.430 (1.281)	0.219 (0.037)
	MCP-BIC	0.857 (0.128)	1.000 (0.000)	1.930 (1.539)	-1.430 (1.281)	0.216 (0.040)
	SCAD-CV	0.858 (0.129)	1.000 (0.000)	3.281 (2.197)	-1.420 (1.288)	1.456 (0.276)
	SCAD-SIC	0.856 (0.130)	1.000 (0.000)	3.289 (2.206)	-1.440 (1.297)	0.220 (0.038)
	SCAD-EBIC	0.858 (0.129)	1.000 (0.000)	3.281 (2.197)	-1.420 (1.288)	0.221 (0.038)
	SCAD-BIC	0.858 (0.129)	1.000 (0.000)	3.281 (2.197)	-1.420 (1.288)	0.221 (0.035)
	LASSO-CV	0.976 (0.047)	0.987 (0.014)	0.737 (0.284)	19.430 (21.005)	1.307 (0.179)
LASSO-SIC	0.969 (0.056)	0.999 (0.001)	1.099 (0.605)	1.100 (1.534)	0.117 (0.029)	
LASSO-EBIC	0.971 (0.050)	0.999 (0.001)	1.068 (0.600)	1.460 (1.772)	0.117 (0.022)	
LASSO-BIC	0.975 (0.048)	0.998 (0.002)	0.947 (0.501)	3.250 (3.255)	0.114 (0.020)	
$p = 2500$	ABESS	0.957 (0.067)	1.000 (0.000)	0.139 (0.177)	0.000 (0.953)	0.523 (0.077)
	SDAR-EBIC	0.957 (0.067)	1.000 (0.000)	0.152 (0.188)	0.090 (0.996)	2.869 (0.538)
	SDAR-SIC	0.957 (0.067)	1.000 (0.000)	0.139 (0.177)	0.000 (0.953)	2.837 (0.533)
	CWM-CV	0.957 (0.067)	1.000 (0.000)	0.207 (0.252)	0.720 (1.138)	7.191 (0.833)
	CWM-SIC	0.957 (0.067)	1.000 (0.000)	0.139 (0.177)	0.000 (0.953)	0.816 (0.102)
	PSIM-CV	0.957 (0.067)	1.000 (0.000)	0.206 (0.252)	0.710 (1.140)	48.601 (5.930)
	PSIM-SIC	0.957 (0.067)	1.000 (0.000)	0.139 (0.177)	0.000 (0.953)	5.248 (0.963)
	MCP-CV	0.836 (0.131)	1.000 (0.000)	2.029 (1.566)	-1.640 (1.314)	2.426 (0.440)
	MCP-SIC	0.835 (0.133)	1.000 (0.000)	2.030 (1.567)	-1.650 (1.329)	0.344 (0.062)
	MCP-EBIC	0.835 (0.133)	1.000 (0.000)	2.030 (1.567)	-1.650 (1.329)	0.342 (0.059)
	MCP-BIC	0.836 (0.131)	1.000 (0.000)	2.029 (1.566)	-1.640 (1.314)	0.342 (0.059)
	SCAD-CV	0.836 (0.131)	1.000 (0.000)	3.563 (2.717)	-1.640 (1.314)	2.363 (0.386)
	SCAD-SIC	0.834 (0.132)	1.000 (0.000)	3.567 (2.716)	-1.660 (1.320)	0.336 (0.061)
	SCAD-EBIC	0.836 (0.131)	1.000 (0.000)	3.563 (2.717)	-1.640 (1.314)	0.333 (0.058)
	SCAD-BIC	0.836 (0.131)	1.000 (0.000)	3.563 (2.717)	-1.640 (1.314)	0.335 (0.055)
	LASSO-CV	0.956 (0.069)	0.991 (0.010)	0.811 (0.487)	23.020 (26.117)	2.019 (0.281)
LASSO-SIC	0.951 (0.069)	1.000 (0.001)	1.212 (0.887)	0.710 (1.526)	0.174 (0.027)	
LASSO-EBIC	0.953 (0.069)	0.999 (0.001)	1.188 (0.877)	0.910 (1.688)	0.177 (0.029)	
LASSO-BIC	0.954 (0.069)	0.999 (0.002)	1.034 (0.709)	3.130 (4.007)	0.177 (0.032)	

Table S3. TPR, TNR, relative errors and sparse level errors of nineteen variable selection methods for constant correlation structure.

	Method	TPR	TNR	ReErr	SLE	Runtime
$p = 500$	ABESS	0.907 (0.091)	0.999 (0.001)	0.669 (0.615)	-0.430 (1.121)	0.098 (0.026)
	SDAR-EBIC	0.900 (0.097)	0.999 (0.001)	0.623 (0.599)	-0.670 (1.092)	0.511 (0.172)
	SDAR-SIC	0.906 (0.092)	0.999 (0.001)	0.674 (0.612)	-0.430 (1.121)	0.507 (0.162)
	CWM-CV	0.912 (0.089)	0.991 (0.009)	1.602 (1.798)	3.500 (4.629)	3.371 (0.765)
	CWM-SIC	0.898 (0.097)	0.999 (0.002)	0.686 (0.657)	-0.420 (1.273)	0.319 (0.086)
	PSIM-CV	0.914 (0.090)	0.993 (0.008)	1.459 (1.653)	2.470 (4.279)	24.206 (5.539)
	PSIM-SIC	0.900 (0.097)	0.999 (0.002)	0.691 (0.638)	-0.430 (1.257)	2.406 (0.614)
	MCP-CV	0.542 (0.227)	1.000 (0.000)	52.418 (61.560)	-4.580 (2.270)	0.316 (0.055)
	MCP-SIC	0.540 (0.227)	1.000 (0.000)	52.419 (61.560)	-4.600 (2.270)	0.042 (0.014)
	MCP-EBIC	0.540 (0.227)	1.000 (0.000)	52.419 (61.560)	-4.600 (2.270)	0.041 (0.016)
	MCP-BIC	0.540 (0.227)	1.000 (0.000)	52.419 (61.560)	-4.600 (2.270)	0.042 (0.015)
	SCAD-CV	0.552 (0.227)	1.000 (0.000)	72.049 (85.035)	-4.480 (2.267)	0.331 (0.066)
	SCAD-SIC	0.542 (0.220)	1.000 (0.000)	72.101 (85.009)	-4.580 (2.198)	0.041 (0.015)
	SCAD-EBIC	0.544 (0.223)	1.000 (0.000)	72.091 (85.016)	-4.560 (2.231)	0.046 (0.017)
	SCAD-BIC	0.550 (0.227)	1.000 (0.000)	72.053 (85.033)	-4.500 (2.267)	0.043 (0.014)
	LASSO-CV	0.847 (0.119)	0.987 (0.023)	11.228 (10.250)	4.910 (11.914)	0.440 (0.123)
LASSO-SIC	0.901 (0.090)	0.991 (0.008)	4.484 (3.158)	3.270 (4.080)	0.098 (0.066)	
LASSO-EBIC	0.905 (0.088)	0.989 (0.009)	4.089 (2.978)	4.350 (4.635)	0.102 (0.073)	
LASSO-BIC	0.915 (0.083)	0.985 (0.009)	3.437 (2.350)	6.550 (4.751)	0.098 (0.070)	
$p = 1500$	ABESS	0.930 (0.088)	1.000 (0.001)	0.768 (0.851)	-0.240 (1.084)	0.266 (0.065)
	SDAR-EBIC	0.929 (0.089)	1.000 (0.001)	0.764 (0.855)	-0.270 (1.081)	1.239 (0.266)
	SDAR-SIC	0.930 (0.088)	1.000 (0.001)	0.768 (0.851)	-0.240 (1.084)	1.229 (0.261)
	CWM-CV	0.933 (0.082)	0.997 (0.004)	1.541 (1.503)	3.420 (5.919)	3.525 (0.727)
	CWM-SIC	0.920 (0.086)	1.000 (0.001)	0.850 (0.878)	-0.120 (1.305)	0.351 (0.067)
	PSIM-CV	0.935 (0.081)	0.997 (0.004)	1.885 (2.363)	3.290 (6.243)	54.868 (12.254)
	PSIM-SIC	0.920 (0.086)	1.000 (0.001)	0.881 (0.902)	-0.080 (1.331)	6.325 (1.944)
	MCP-CV	0.540 (0.204)	1.000 (0.000)	53.214 (63.374)	-4.600 (2.045)	1.103 (0.232)
	MCP-SIC	0.540 (0.204)	1.000 (0.000)	53.214 (63.374)	-4.600 (2.045)	0.161 (0.043)
	MCP-EBIC	0.540 (0.204)	1.000 (0.000)	53.214 (63.374)	-4.600 (2.045)	0.162 (0.039)
	MCP-BIC	0.540 (0.204)	1.000 (0.000)	53.214 (63.374)	-4.600 (2.045)	0.163 (0.042)
	SCAD-CV	0.540 (0.208)	1.000 (0.000)	81.637 (87.852)	-4.600 (2.084)	1.107 (0.208)
	SCAD-SIC	0.533 (0.207)	1.000 (0.000)	81.785 (87.959)	-4.670 (2.070)	0.160 (0.041)
	SCAD-EBIC	0.533 (0.207)	1.000 (0.000)	81.785 (87.959)	-4.670 (2.070)	0.165 (0.036)
	SCAD-BIC	0.540 (0.208)	1.000 (0.000)	81.637 (87.853)	-4.600 (2.084)	0.163 (0.041)
	LASSO-CV	0.861 (0.136)	0.994 (0.011)	13.456 (9.583)	7.230 (16.650)	1.265 (0.306)
LASSO-SIC	0.852 (0.138)	0.999 (0.002)	14.339 (8.904)	0.140 (3.260)	0.112 (0.033)	
LASSO-EBIC	0.855 (0.137)	0.999 (0.002)	14.135 (9.017)	0.530 (3.878)	0.113 (0.034)	
LASSO-BIC	0.859 (0.136)	0.998 (0.003)	13.747 (9.322)	1.620 (5.065)	0.112 (0.032)	
$p = 2500$	ABESS	0.909 (0.087)	1.000 (0.000)	0.885 (0.871)	-0.190 (1.253)	0.460 (0.093)
	SDAR-EBIC	0.909 (0.087)	1.000 (0.000)	0.958 (1.006)	-0.080 (1.316)	1.917 (0.334)
	SDAR-SIC	0.909 (0.087)	1.000 (0.000)	0.887 (0.870)	-0.180 (1.258)	1.908 (0.335)
	CWM-CV	0.907 (0.088)	0.998 (0.003)	1.502 (1.325)	3.580 (7.087)	4.492 (0.785)
	CWM-SIC	0.897 (0.094)	1.000 (0.000)	0.930 (0.937)	-0.240 (1.379)	0.461 (0.088)
	PSIM-CV	0.908 (0.085)	0.999 (0.002)	1.402 (1.430)	1.620 (4.456)	92.834 (21.810)
	PSIM-SIC	0.904 (0.091)	1.000 (0.000)	0.943 (0.950)	-0.090 (1.386)	11.009 (2.951)
	MCP-CV	0.512 (0.199)	1.000 (0.000)	55.665 (58.164)	-4.880 (1.991)	2.172 (0.429)
	MCP-SIC	0.512 (0.199)	1.000 (0.000)	55.665 (58.164)	-4.880 (1.991)	0.291 (0.061)
	MCP-EBIC	0.512 (0.199)	1.000 (0.000)	55.665 (58.164)	-4.880 (1.991)	0.289 (0.074)
	MCP-BIC	0.512 (0.199)	1.000 (0.000)	55.665 (58.164)	-4.880 (1.991)	0.282 (0.064)
	SCAD-CV	0.512 (0.195)	1.000 (0.000)	84.011 (87.386)	-4.880 (1.945)	2.203 (0.427)
	SCAD-SIC	0.501 (0.196)	1.000 (0.000)	84.259 (87.703)	-4.990 (1.962)	0.292 (0.075)
	SCAD-EBIC	0.504 (0.194)	1.000 (0.000)	84.074 (87.383)	-4.960 (1.943)	0.290 (0.063)
	SCAD-BIC	0.507 (0.195)	1.000 (0.000)	84.037 (87.386)	-4.930 (1.950)	0.282 (0.062)
	LASSO-CV	0.851 (0.134)	0.996 (0.007)	12.895 (8.866)	7.900 (17.620)	2.310 (0.584)
LASSO-SIC	0.843 (0.131)	0.999 (0.002)	14.099 (8.597)	0.900 (4.800)	0.202 (0.057)	
LASSO-EBIC	0.845 (0.132)	0.999 (0.002)	13.649 (8.647)	1.600 (6.145)	0.206 (0.062)	
LASSO-BIC	0.848 (0.134)	0.998 (0.003)	13.448 (8.778)	2.300 (6.927)	0.206 (0.062)	

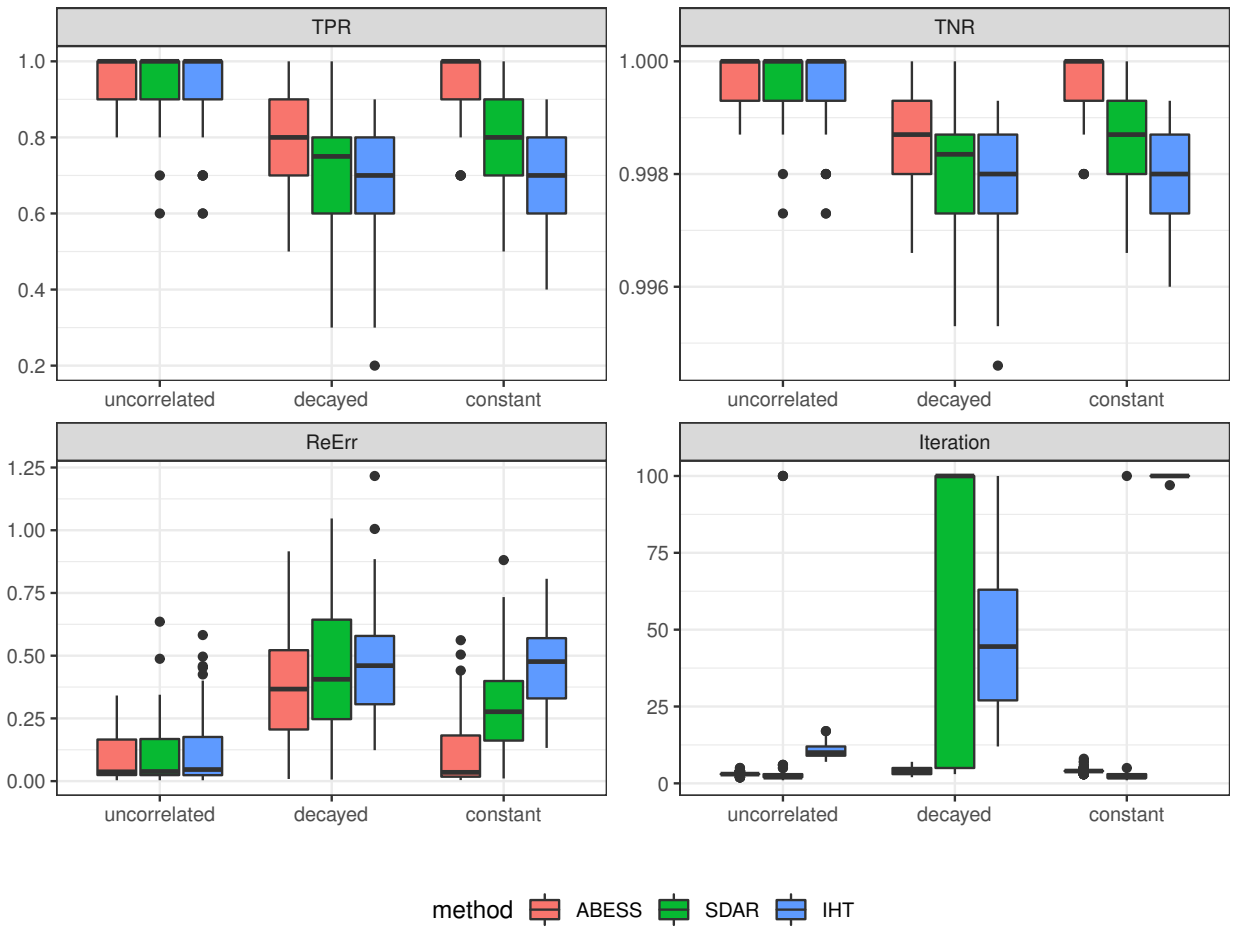


Fig. S1. TNR, TPR, relative error, and the number of iterations in the fixed support size best subset selection problem.

71 can efficiently speed up our algorithm. We conduct numerical simulation studies to explore the utility of the GS strategy. The
 72 simulation settings are completely adopted from Section Simulation. The numerical results are displayed in Figure S2. As we
 73 can see in Figure S2, by combining golden section search and SIC (or CV), we can accelerate the best sparsity level selection
 procedure while maintaining the sparsity level selection performance.

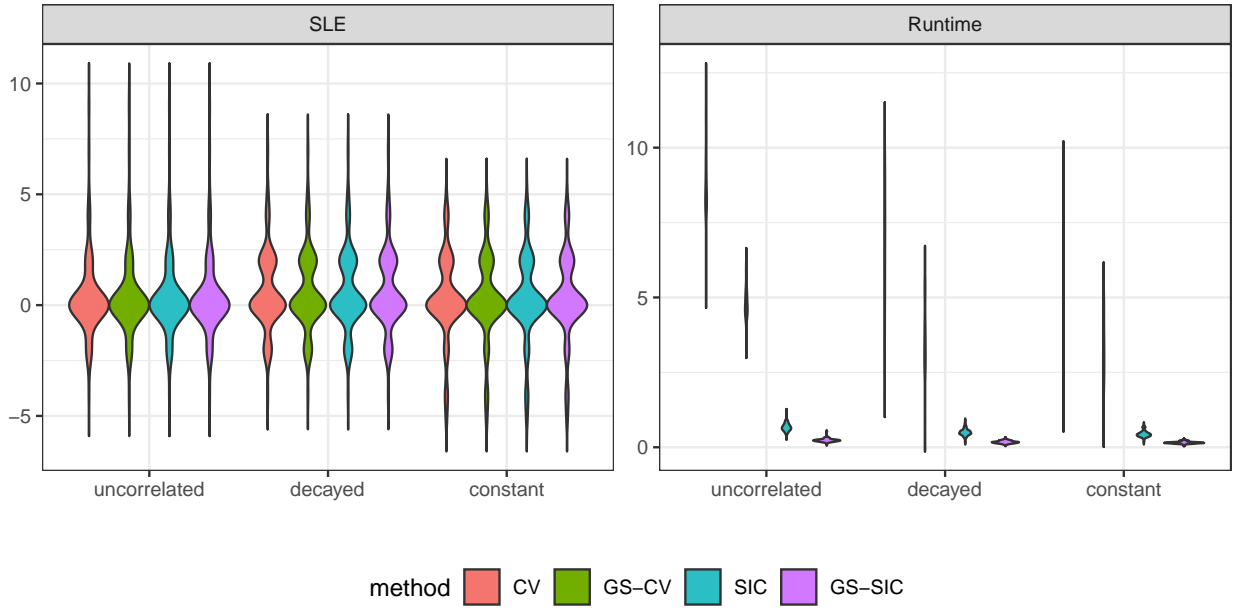


Fig. S2. Sparsity level error and runtime comparison in the sparsity level selection problem. The methods with prefix “GS” is golden section based the sparsity level selection strategy.

74

75 Proofs of Main results

76 Proof of Lemma 1.

proof 1 *Let*

$$\mathcal{A}_1 = \hat{\mathcal{A}} \cap \mathcal{A}^*, \mathcal{A}_2 = \hat{\mathcal{A}} \cap \mathcal{I}^*$$

$$\mathcal{I}_1 = \hat{\mathcal{I}} \cap \mathcal{A}^*, \mathcal{I}_2 = \hat{\mathcal{I}} \cap \mathcal{I}^*$$

77 *We assume $\mathcal{I}_1 \neq \emptyset$ and show that it will lead to a contradiction.*

Let $k = |\mathcal{I}_1|$. Denote the splicing set in the active and inactive sets, respectively, as

$$\hat{\mathcal{A}}_k = \{j \in \hat{\mathcal{A}} : \sum_{i \in \hat{\mathcal{A}}} \mathbb{I}(|\hat{\beta}_j| \geq |\hat{\beta}_i|) \leq k\},$$

$$\hat{\mathcal{I}}_k = \{j \in \hat{\mathcal{I}} : \sum_{i \in \hat{\mathcal{I}}} \mathbb{I}(|\hat{d}_j| \leq |\hat{d}_i|) \leq k\}.$$

And denote

$$\mathcal{A}_{11} = \mathcal{A}_1 \cap (\hat{\mathcal{A}}_k)^c, \mathcal{A}_{12} = \mathcal{A}_1 \cap \hat{\mathcal{A}}_k,$$

$$\mathcal{A}_{21} = \mathcal{A}_2 \cap (\hat{\mathcal{A}}_k)^c, \mathcal{A}_{22} = \mathcal{A}_2 \cap \hat{\mathcal{A}}_k,$$

and

$$\mathcal{I}_{11} = \mathcal{I}_1 \cap \hat{\mathcal{I}}_k, \mathcal{I}_{12} = \mathcal{I}_1 \cap (\hat{\mathcal{I}}_k)^c,$$

$$\mathcal{I}_{21} = \mathcal{I}_2 \cap \hat{\mathcal{I}}_k, \mathcal{I}_{22} = \mathcal{I}_2 \cap (\hat{\mathcal{I}}_k)^c.$$

Consider the following four cases:

- (1) $\mathcal{I}_{12} \neq \emptyset, \mathcal{A}_{12} \neq \emptyset,$ (2) $\mathcal{I}_{12} \neq \emptyset, \mathcal{A}_{12} = \emptyset,$
- (3) $\mathcal{I}_{12} = \emptyset, \mathcal{A}_{12} \neq \emptyset,$ (4) $\mathcal{I}_{12} = \emptyset, \mathcal{A}_{12} = \emptyset.$

78 We provide the details for the first case as the other cases follow similarly.

Let $\mathbf{H}_A = \mathbf{X}_A (\mathbf{X}'_A \mathbf{X}_A)^{-1} \mathbf{X}'_A$, for any index set $A \subseteq \{1, \dots, p\}$. The estimator $(\hat{\beta}, \hat{\mathbf{d}})$ can be expressed as

$$\begin{aligned}\hat{\beta}_{\mathcal{A}_1} &= (\mathbf{X}'_{\mathcal{A}_1} (\mathbf{I} - \mathbf{H}_{\mathcal{A}_2}) \mathbf{X}_{\mathcal{A}_1})^{-1} \mathbf{X}'_{\mathcal{A}_1} (\mathbf{I} - \mathbf{H}_{\mathcal{A}_2}) \mathbf{y}, \\ \hat{\beta}_{\mathcal{A}_2} &= (\mathbf{X}'_{\mathcal{A}_2} (\mathbf{I} - \mathbf{H}_{\mathcal{A}_1}) \mathbf{X}_{\mathcal{A}_2})^{-1} \mathbf{X}'_{\mathcal{A}_2} (\mathbf{I} - \mathbf{H}_{\mathcal{A}_1}) \mathbf{y}, \\ \hat{\mathbf{d}}_{\mathcal{I}_1} &= \mathbf{X}'_{\mathcal{I}_1} (\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}}) \mathbf{y} / n, \\ \hat{\mathbf{d}}_{\mathcal{I}_2} &= \mathbf{X}'_{\mathcal{I}_2} (\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}}) \mathbf{y} / n.\end{aligned}$$

79 First of all, Assume that events $\{c_-(s) \|\beta_{\mathcal{I}_2}^*\|_2 \leq 2(1 + \Delta) \left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)} \right) \|\beta_{\mathcal{I}_1}^*\|_2\}$ and $\{\|\beta_{\mathcal{A}_2}^*\|_2 \leq 2(1 + \Delta) \frac{\theta_{s,s}}{c_-(s)} \|\beta_{\mathcal{I}_1}^*\|_2\}$
80 hold. We will show that these two events hold with probability $1 - \frac{1}{3}\gamma_1 - \frac{1}{3}\gamma_2$ later.

81 Now, we splicing $\hat{\mathcal{A}}_k = \mathcal{A}_{12} \cup \mathcal{A}_{22}$ and $\hat{\mathcal{I}}_k = \mathcal{I}_{11} \cup \mathcal{I}_{21}$, and then the new active set is $\tilde{\mathcal{A}} = (\hat{\mathcal{A}} \setminus \hat{\mathcal{A}}_k) \cup \hat{\mathcal{I}}_k$ and the inactive set
82 is $\tilde{\mathcal{I}} = (\tilde{\mathcal{A}})^c$. Let $\tilde{\beta} = \arg \min_{\beta_{\tilde{\mathcal{I}}=0} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2}$. The loss function of $\tilde{\beta}$ is

$$\begin{aligned}2n\mathcal{L}(\tilde{\beta}) &= \|\mathbf{y} - \mathbf{X}\tilde{\beta}\|_2^2 = \mathbf{y}' (\mathbf{I} - \mathbf{H}_{\tilde{\mathcal{A}}}) \mathbf{y} \\ &= (\mathbf{X}_{\mathcal{A}_{12} \cup \mathcal{I}_{12}} \beta_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}^* + \epsilon)' (\mathbf{I} - \mathbf{H}_{\tilde{\mathcal{A}}}) (\mathbf{X}_{\mathcal{A}_{12} \cup \mathcal{I}_{12}} \beta_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}^* + \epsilon) \\ &\leq nc_+(s) (\|\beta_{\mathcal{A}_{12}}^*\|_2^2 + \|\beta_{\mathcal{I}_{12}}^*\|_2^2) + \epsilon' \epsilon + \\ &\quad 2|\epsilon' (\mathbf{I} - \mathbf{H}_{\tilde{\mathcal{A}}}) \mathbf{X}_{\mathcal{A}_{12} \cup \mathcal{I}_{12}} \beta_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}^*| + |\epsilon' \mathbf{H}_{\tilde{\mathcal{A}}} \epsilon| \\ &\leq nc_+(s) \left[\left(2(1 + \Delta) \frac{\theta_{s,s}}{c_-(s)} \right)^2 + \left(2(1 + \Delta) \frac{\theta_{s,s}}{c_-(s)} \left(1 + \frac{\theta_{s,s}}{c_-(s)} \right) \right)^2 \right] \|\beta_{\mathcal{I}_1}^*\|_2^2 + \\ &\quad \epsilon' \epsilon + 2|\epsilon' (\mathbf{I} - \mathbf{H}_{\tilde{\mathcal{A}}}) \mathbf{X}_{\mathcal{A}_{12} \cup \mathcal{I}_{12}} \beta_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}^*| + |\epsilon' \mathbf{H}_{\tilde{\mathcal{A}}} \epsilon| \\ &\leq 8nc_+(s) \left(\left(1 + \Delta \right) \frac{\theta_{s,s}}{c_-(s)} \left(1 + \frac{\theta_{s,s}}{c_-(s)} \right) \right)^2 \|\beta_{\mathcal{I}_1}^*\|_2^2 + \epsilon' \epsilon + f_1(\epsilon),\end{aligned}\tag{1}$$

84 where $f_1(\epsilon) = 2|\epsilon' (\mathbf{I} - \mathbf{H}_{\tilde{\mathcal{A}}}) \mathbf{X}_{\mathcal{A}_{12} \cup \mathcal{I}_{12}} \beta_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}^*| + |\epsilon' \mathbf{H}_{\tilde{\mathcal{A}}} \epsilon|$. In addition, the loss function of $\hat{\beta}$ is

$$\begin{aligned}2n\mathcal{L}(\hat{\beta}) &= \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 = \mathbf{y}' (\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}}) \mathbf{y} \\ &= (\mathbf{X}_{\mathcal{I}_1} \beta_{\mathcal{I}_1}^* + \epsilon)' (\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}}) (\mathbf{X}_{\mathcal{I}_1} \beta_{\mathcal{I}_1}^* + \epsilon) \\ &\geq n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1}^*\|_2^2 + \epsilon' \epsilon - 2|\epsilon' (\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}}) \mathbf{X}_{\mathcal{I}_1} \beta_{\mathcal{I}_1}^*| - |\epsilon' \mathbf{H}_{\hat{\mathcal{A}}} \epsilon| \\ &= n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1}^*\|_2^2 + \epsilon' \epsilon - f_2(\epsilon),\end{aligned}\tag{2}$$

where $f_2(\epsilon) = 2|\epsilon' (\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}}) \mathbf{X}_{\mathcal{I}_1} \beta_{\mathcal{I}_1}^*| + |\epsilon' \mathbf{H}_{\hat{\mathcal{A}}} \epsilon|$. The conditions of this lemma assure that

$$(1 - \Delta) n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1}^*\|_2^2 > 8nc_+(s) \left(\left(1 + \Delta \right) \frac{\theta_{s,s}}{c_-(s)} \left(1 + \frac{\theta_{s,s}}{c_-(s)} \right) \right)^2 \|\beta_{\mathcal{I}_1}^*\|_2^2.$$

Let $\mathbf{u}_j = e_j' (\mathbf{X}'_{\mathcal{A}} \mathbf{X}_{\mathcal{A}})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{A}}$, where the j th element of \mathbf{e}_j is 1 and otherwise is 0. Note that $\|\mathbf{u}_j\|_2^2 \leq 1$, we have

$$\begin{aligned}\mathbb{P} \left(\|(\mathbf{X}'_{\mathcal{A}} \mathbf{X}_{\mathcal{A}})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{A}} \epsilon\|_2 \geq t \right) &\leq \sum_{j \in \mathcal{A}} \mathbb{P} \left(\|e_j' (\mathbf{X}'_{\mathcal{A}} \mathbf{X}_{\mathcal{A}})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{A}} \epsilon\|_2 \geq \frac{t}{\sqrt{|\mathcal{A}|}} \right) \\ &\leq \sum_{j \in \mathcal{A}} \mathbb{P} \left(\|\mathbf{u}_j \epsilon\|_2 \geq \frac{t}{\sqrt{|\mathcal{A}|}} \right) \\ &\leq 2p \exp \left\{ -\frac{t^2}{\sigma^2 |\mathcal{A}|} \right\}.\end{aligned}$$

Thus

$$\begin{aligned}\mathbb{P} \left(\|\mathbf{H}_{\tilde{\mathcal{A}}} \epsilon\|_2 \geq \sqrt{\frac{\Delta}{4} n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})} \|\beta_{\mathcal{I}_1}^*\|_2 \right) &\leq \frac{\gamma_3}{6}, \\ \mathbb{P} \left(\|\mathbf{H}_{\hat{\mathcal{A}}} \epsilon\|_2 \geq \sqrt{\frac{\Delta}{4} n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})} \|\beta_{\mathcal{I}_1}^*\|_2 \right) &\leq \frac{\gamma_3}{6}\end{aligned}$$

where $\gamma_3 = 12 \exp\{\log p - \frac{K_{s,3}nb^*}{s}\}$, $K_{s,3} = \frac{\Delta(c_-(s)^2 - \theta_{s,s}^2)}{4c_-(s)\sigma^2}$. Similarly, we can show that,

$$\begin{aligned} \mathbb{P}\left(2|\epsilon'(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\mathbf{X}_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}\beta_{\mathcal{A}_{12} \cup \mathcal{I}_{12}}^*| \geq \frac{\Delta}{4}n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})\|\beta_{\mathcal{I}_1}^*\|_2^2\right) &\leq \frac{\gamma_4}{2}, \\ \mathbb{P}\left(2|\epsilon'(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^*| \geq \frac{\Delta}{4}n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})\|\beta_{\mathcal{I}_1}^*\|_2^2\right) &\leq \frac{\gamma_4}{2}, \end{aligned}$$

where $\gamma_4 = 4 \exp\{\log p - \frac{K_{s,4}nb^*}{s^*}\}$, $K_{s,4} = \min\left\{\frac{(\Delta(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})/8)^2}{c_+(s)\sigma^2} / \left(4(1 + \Delta)\frac{\theta_{s,s}}{c_-(s)}\left(1 + \frac{\theta_{s,s}}{c_-(s)}\right)\right)^2, \frac{(\Delta(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})/8)^2}{c_+(s)\sigma^2}\right\}$. Thus,

$$\mathbb{P}\left(f_1(\epsilon) + f_2(\epsilon) \geq \Delta n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})\|\beta_{\mathcal{I}_1}^*\|_2^2\right) \leq \frac{1}{3}\gamma_3 + \gamma_4.$$

Now, it follows from Conditions (4) and (6) that

$$\begin{aligned} \mathcal{L}(\hat{\beta}) - \mathcal{L}(\tilde{\beta}) &\geq \frac{(1 - \delta_s)(1 - \Delta)(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})}{2}\|\beta_{\mathcal{I}_1}^*\|_2^2 \\ &\geq \frac{(1 - \delta_s)(1 - \Delta)(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})}{2}b^* \\ &> \tau_s. \end{aligned}$$

Consequently,

$$\mathbb{P}\left(\mathcal{L}(\hat{\beta}) - \mathcal{L}(\tilde{\beta}) > \tau_s\right) \geq 1 - \frac{1}{3}\gamma_1 - \frac{1}{3}\gamma_2 - \frac{1}{3}\gamma_3 - \gamma_4,$$

which leads to a contradiction with $\mathcal{I}_1 \neq \emptyset$. Therefore,

$$\mathbb{P}(\hat{\mathcal{A}} \supseteq \mathcal{A}^*) \geq 1 - \gamma(s, n, p, b^*),$$

where $\gamma(s, n, p, b^*) = 16 \exp\{\log p - \frac{K_s nb^*}{s}\}$ and $K_s = \min\{K_{s,1}, K_{s,2}, K_{s,3}, K_{s,4}\}$. By Condition (6),

$$\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\mathcal{A}} \supseteq \mathcal{A}^*) = 1.$$

Especially, if $s = s^*$, we can show that

$$\mathbb{P}(\hat{\mathcal{A}} = \mathcal{A}^*) = 1.$$

It remains to show that

$$\mathbb{P}\left(c_-(s)\|\beta_{\mathcal{I}_{12}}^*\|_2 \leq 2(1 + \Delta)\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2\right) \geq 1 - \frac{1}{3}\gamma_1. \quad [3]$$

and

$$\mathbb{P}\left(\|\beta_{\mathcal{A}_{12}}^*\|_2 \leq 2(1 + \Delta)\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right) \geq 1 - \frac{1}{3}\gamma_2. \quad [4]$$

On one hand, it follows from $|\hat{\mathcal{I}}_k| = |\mathcal{I}_{11}| + |\mathcal{I}_{21}| = |\mathcal{I}_1|$, $|\mathcal{I}_{12}| = |\mathcal{I}_{21}|$ and the definition of $\hat{\mathcal{I}}_k$ that

$$\min_{j \in \mathcal{I}_{21}} |\hat{d}_j| \geq \max_{j \in \mathcal{I}_{12}} |\hat{d}_j|.$$

Note that,

$$\begin{aligned} n\|\hat{d}_{\mathcal{I}_{12}}\|_2 &= \|\mathbf{X}'_{\mathcal{I}_{12}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\mathbf{y}\|_2 \\ &= \|\mathbf{X}'_{\mathcal{I}_{12}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})(\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^* + \epsilon)\|_2 \\ &\geq \|\mathbf{X}'_{\mathcal{I}_{12}}\mathbf{X}_{\mathcal{I}_{12}}\beta_{\mathcal{I}_{12}}^*\|_2 - \|\mathbf{X}'_{\mathcal{I}_{12}}\mathbf{X}_{\mathcal{I}_{11}}\beta_{\mathcal{I}_{11}}^*\|_2 - \\ &\quad \|\mathbf{X}'_{\mathcal{I}_{12}}\mathbf{H}_{\hat{\mathcal{A}}}\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^*\|_2 - \|\mathbf{X}'_{\mathcal{I}_{12}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2 \\ &\geq nc_-(s)\|\beta_{\mathcal{I}_{12}}^*\|_2 - n\theta_{s,s}\|\beta_{\mathcal{I}_{11}}^*\|_2 - n\frac{\theta_{s,s}^2}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2 - \|\mathbf{X}'_{\mathcal{I}_{12}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2 \end{aligned} \quad [5]$$

and

$$\begin{aligned} n\|\hat{d}_{\mathcal{I}_{21}}\|_2 &= \|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\mathbf{y}\|_2 \\ &= \|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})(\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^* + \epsilon)\|_2 \\ &\leq \|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^*\|_2 + \|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2 \\ &\leq n\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2 + \|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2. \end{aligned} \quad [6]$$

Since $|\mathcal{I}_{12}| = |\mathcal{I}_{21}|$, $\|\hat{\mathbf{d}}_{\mathcal{I}_{21}}\|_2 \geq \|\hat{\mathbf{d}}_{\mathcal{I}_{12}}\|_2$, combined with (5) and (6), we have

$$\begin{aligned} & 2n\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2 + \|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2 \\ & \geq nc_-(s)\|\beta_{\mathcal{I}_{12}}^*\|_2 - \|\mathbf{X}'_{\mathcal{I}_{12}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2. \end{aligned}$$

Let $\mathbf{h}'_j = \mathbf{X}'_j(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})$,

$$\|\mathbf{h}_j\|_2^2 = |\mathbf{X}'_j(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\mathbf{X}_j| \leq nc_+(s).$$

By Condition (1) and $|\mathcal{I}_1| \geq |\mathcal{I}_{21}|$, for some $\Delta > 0$, we have

$$\begin{aligned} & \mathbb{P}\left(\|\mathbf{X}'_{\mathcal{I}_{21}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2 > n\Delta\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ & \leq \sum_{j \in \mathcal{I}_{21}} \mathbb{P}\left(|\mathbf{h}'_j \epsilon| > \frac{n\Delta\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2}{\sqrt{|\mathcal{I}_{21}|}}\right) \\ & \leq 2p \exp\left\{-\frac{1}{nc_+(s)\sigma^2} \frac{n^2\Delta^2\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)^2}{|\mathcal{I}_{21}|} \|\beta_{\mathcal{I}_1}^*\|_2^2\right\} \\ & \leq 2p \exp\left\{-\frac{n^2\Delta^2\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)^2}{nc_+(s)\sigma^2} (\min_{j \in \mathcal{A}^*} |\beta_j^*|)^2\right\} \\ & = \frac{1}{6}\gamma_1, \end{aligned}$$

where $\gamma_1 = 12 \exp\{\log p - K_{s,1}nb^*\}$ and $K_{s,1} = \frac{\Delta^2\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)^2}{c_+(s)\sigma^2}$. Similarly,

$$\begin{aligned} & \mathbb{P}\left(\|\mathbf{X}'_{\mathcal{I}_{12}}(\mathbf{I} - \mathbf{H}_{\hat{\mathcal{A}}})\epsilon\|_2 > n\Delta\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ & \leq 2p \exp\left\{-\frac{1}{2nc_+(s)\sigma^2} \frac{n^2\Delta^2\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)^2}{|\mathcal{I}_{12}|} \|\beta_{\mathcal{I}_1}^*\|_2^2\right\} \\ & \leq \frac{1}{6}\gamma_1. \end{aligned}$$

Consequently,

$$\mathbb{P}\left(c_-(s)\|\beta_{\mathcal{I}_{12}}^*\|_2 \leq 2(1 + \Delta)\left(\theta_{s,s} + \frac{\theta_{s,s}^2}{c_-(s)}\right)\|\beta_{\mathcal{I}_1}^*\|_2\right) \geq 1 - \frac{1}{3}\gamma_1.$$

On the other hand,

$$\begin{aligned} \hat{\beta}_{\mathcal{A}_1} &= \beta_{\mathcal{A}_1}^* + (\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1} \mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})(\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^* + \epsilon), \\ \hat{\beta}_{\mathcal{A}_2} &= 0 + (\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\mathbf{X}_{\mathcal{A}_2})^{-1} \mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})(\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^* + \epsilon). \end{aligned}$$

Then $\hat{\mathcal{A}}_k = \mathcal{A}_{12} \cup \mathcal{A}_{22}$ and $|\mathcal{A}_{12}| + |\mathcal{A}_{22}| = |\mathcal{I}_1| \leq |\mathcal{A}_2| = |\mathcal{A}_{21}| + |\mathcal{A}_{22}|$, and we have

$$\max_{j \in \mathcal{A}_{12} \cup \mathcal{A}_{22}} |\hat{\beta}_j| \leq \min_{j \in \mathcal{A}_{11} \cup \mathcal{A}_{21}} |\hat{\beta}_j|.$$

Considering the case when $\mathcal{A}_{12} \neq \emptyset$, we have

$$\frac{1}{\sqrt{|\mathcal{A}_{12}|}} \|\hat{\beta}_{\mathcal{A}_{12}}\|_2 \leq \frac{1}{\sqrt{|\mathcal{A}_{21}|}} \|\hat{\beta}_{\mathcal{A}_{21}}\|_2.$$

94 Denote $\mathbf{E}_{\mathcal{A}_{12}}$ as a $|\mathcal{A}_{12}| \times |\mathcal{A}_1|$ matrix and its j th row is a $|\mathcal{A}_1|$ -dimensional vector \mathbf{e}_j , where the j th element of \mathbf{e}_j is 1 and
95 otherwise 0. $\mathbf{E}_{\mathcal{A}_{21}}$ can be defined analogously. Then

$$\begin{aligned} \|\hat{\beta}_{\mathcal{A}_{12}}\|_2 & \geq \|\beta_{\mathcal{A}_{12}}^*\|_2 - \|\mathbf{E}_{\mathcal{A}_{12}}(\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1} \mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})(\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^* + \epsilon)\|_2 \\ & \geq \|\beta_{\mathcal{A}_{12}}^*\|_2 - [nc_-(s)]^{-1} n\theta_{s,s} \|\beta_{\mathcal{I}_1}^*\|_2 - \|\mathbf{E}_{\mathcal{A}_{12}}(\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1} \mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\epsilon\|_2, \end{aligned} \quad [7]$$

97 and

$$\begin{aligned} \|\hat{\beta}_{\mathcal{A}_{21}}\|_2 &\leq \|\mathbf{E}_{\mathcal{A}_{21}}(\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\mathbf{X}_{\mathcal{A}_2})^{-1}\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})(\mathbf{X}_{\mathcal{I}_1}\beta_{\mathcal{I}_1}^* + \epsilon)\|_2 \\ &\leq [nc_-(s)]^{-1}n\theta_{s,s}\|\beta_{\mathcal{I}_1}^*\|_2 + \|\mathbf{E}_{\mathcal{A}_{21}}(\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\mathbf{X}_{\mathcal{A}_2})^{-1}\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\epsilon\|_2. \end{aligned} \quad [8]$$

Combining (7) and (8), we have

$$\begin{aligned} \sqrt{\frac{|\mathcal{A}_{21}|}{|\mathcal{A}_{12}|}}\|\beta_{\mathcal{A}_{12}}^*\|_2 &\leq \left(1 + \sqrt{\frac{|\mathcal{A}_{21}|}{|\mathcal{A}_{12}|}}\right) \frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2 + \\ &\quad \sqrt{\frac{|\mathcal{A}_{21}|}{|\mathcal{A}_{12}|}}\|\mathbf{E}_{\mathcal{A}_{12}}(\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1}\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\epsilon\|_2 + \\ &\quad \|\mathbf{E}_{\mathcal{A}_{21}}(\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\mathbf{X}_{\mathcal{A}_2})^{-1}\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\epsilon\|_2. \end{aligned}$$

Denote $\mathbf{h}'_j = \mathbf{e}'_j(\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1}\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})$,

$$\|\mathbf{h}_j\|_2^2 = \|\mathbf{e}'_j(\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1}\mathbf{e}_j\|_2^2 \leq [nc_-(s)]^{-1}.$$

It follows that, for some constant $\Delta > 0$,

$$\begin{aligned} &\mathbb{P}\left(\sqrt{\frac{|\mathcal{A}_{21}|}{|\mathcal{A}_{12}|}}\|\mathbf{E}_{\mathcal{A}_{12}}(\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\mathbf{X}_{\mathcal{A}_1})^{-1}\mathbf{X}'_{\mathcal{A}_1}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_2})\epsilon\|_2 > \Delta\sqrt{\frac{|\mathcal{A}_{21}|}{|\mathcal{A}_{12}|}}\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ &\leq \sum_{j \in \mathcal{A}_{12}} \mathbb{P}\left(|\mathbf{h}'_j\epsilon| > \frac{\Delta}{\sqrt{|\mathcal{A}_{12}|}}\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ &\leq 2p \exp\left\{-\frac{nc_-(s)}{\sigma^2}\left(\frac{1}{\sqrt{|\mathcal{A}_{12}|}}\frac{\Delta\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right)^2\right\} \\ &\leq 2p \exp\left\{-\frac{nc_-(s)}{\sigma^2}\left(\sqrt{\frac{|\mathcal{I}_1|}{|\mathcal{A}_{12}|}}\frac{\Delta\theta_{s,s}}{c_-(s)}\min_{j \in \mathcal{A}^*}|\beta_j^*|\right)^2\right\} \\ &= \frac{1}{6}\gamma_2, \end{aligned}$$

where $\gamma_2 = 12 \exp\{\log p - K_{s,2}nb^*\}$ and $K_{s,2} = \frac{(\Delta\theta_{s,s})^2}{\sigma^2c_-(s)}$. Similarly, we have

$$\begin{aligned} &\mathbb{P}\left(\|\mathbf{E}_{\mathcal{A}_{21}}(\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\mathbf{X}_{\mathcal{A}_2})^{-1}\mathbf{X}'_{\mathcal{A}_2}(\mathbf{I} - \mathbf{H}_{\mathcal{A}_1})\epsilon\|_2 > \sqrt{\frac{|\mathcal{A}_{21}|}{|\mathcal{A}_{12}|}}\Delta\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ &\leq 2p \exp\left\{-\frac{nc_-(s)}{\sigma^2}\left(\frac{\Delta}{\sqrt{|\mathcal{A}_{12}|}}\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right)^2\right\} \\ &\leq \frac{1}{6}\gamma_2. \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathbb{P}\left(\|\beta_{\mathcal{A}_{12}}^*\|_2 \leq 2(1 + \Delta)\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ &\geq \mathbb{P}\left(\|\beta_{\mathcal{A}_{12}}^*\|_2 \leq \left(1 + 2\Delta + \frac{\sqrt{|\mathcal{A}_{12}|}}{\sqrt{|\mathcal{A}_{21}|}}\right)\frac{\theta_{s,s}}{c_-(s)}\|\beta_{\mathcal{I}_1}^*\|_2\right) \\ &\geq 1 - \frac{1}{3}\gamma_2. \end{aligned}$$

99 Proof of Theorem 3.

100 **proof 2** It follows from the proof of lemma 1, and inequalities (1) and (2) that

$$\begin{aligned} 2n\mathcal{L}_n(\beta^{m+1}) - 2n\mathcal{L}_n(\beta^*) &\leq 2n\mathcal{L}_n(\tilde{\beta}) - 2n\mathcal{L}_n(\beta^*) \\ &\leq 8nc_+(s)\left((1 + \Delta)\frac{\theta_{s,s}}{c_-(s)}\left(1 + \frac{\theta_{s,s}}{c_-(s)}\right)\right)^2\|\beta_{\mathcal{I}_1^m}^*\|_2^2 + f_1(\epsilon) \\ &= \delta_s(1 - \Delta)n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)})\|\beta_{\mathcal{I}_1^m}^*\|_2^2 + f_1(\epsilon) \\ &\leq \delta_s(2n\mathcal{L}_n(\beta^m) - 2n\mathcal{L}_n(\beta^*)). \end{aligned} \quad [9]$$

101

102 With probability $1 - \gamma(s, n, p, b^*)$. Letting $\mathcal{A}^0 = \emptyset$ and using (9) repeatedly, we have

$$103 \quad 2n\mathcal{L}_n(\beta^m) - 2n\mathcal{L}_n(\beta^*) \leq \delta_s^m (2n\mathcal{L}_n(\beta^0) - 2n\mathcal{L}_n(\beta^*)). \quad [10]$$

104 Now we prove part (i) of the theorem. If $m > \log_{\frac{1}{\delta_s}} \frac{\|\mathbf{y}\|_2^2}{n(1-\Delta) \left(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)^2} \right) b^*} > \log_{\frac{1}{\delta_s}} \left[\frac{|2n\mathcal{L}_n(\beta^0) - 2n\mathcal{L}_n(\beta^*)|}{n(1-\Delta) \left(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)^2} \right) b^*} \right]$,

$$105 \quad 2n\mathcal{L}_n(\beta^m) - 2n\mathcal{L}_n(\beta^*) \leq \delta_s^m (2n\mathcal{L}_n(\beta^0) - 2n\mathcal{L}_n(\beta^*)) < n(1-\Delta) \left(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)^2} \right) b^*. \quad [11]$$

106 Assume $\mathcal{A}^m \not\supseteq \mathcal{A}^*$. ,

$$\begin{aligned} 2n\mathcal{L}_n(\beta^m) - 2n\mathcal{L}_n(\beta^*) &= \|\mathbf{y} - \mathbf{X}\beta^m\|_2^2 - \epsilon'\epsilon = \|(\mathbf{I} - \mathbf{H}_{\mathcal{A}^m})\mathbf{y}\|_2^2 - \epsilon'\epsilon \\ &= \|(\mathbf{I} - \mathbf{H}_{\mathcal{A}^m})(\mathbf{X}_{\mathcal{I}^m}\beta_{\mathcal{I}^m}^m + \epsilon)\|_2^2 - \epsilon'\epsilon \\ &\geq n(1-\Delta) \left(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)^2} \right) b^*, \end{aligned} \quad [12]$$

108 where $\mathbf{H}_{\mathcal{A}^m} = \mathbf{X}_{\mathcal{A}^m}(\mathbf{X}_{\mathcal{A}^m}^\top \mathbf{X}_{\mathcal{A}^m})^{-1} \mathbf{X}_{\mathcal{A}^m}^\top$, $b^* = \min_{j \in \mathcal{A}^*} (\beta_j^*)^2$. Combining (11) and (12) leads to a contradiction. Therefore

109 $\mathcal{A}^m \supseteq \mathcal{A}^*$.

To prove part (ii) of the theorem, note that

$$\begin{aligned} \|\beta^m - \beta^*\|_2 &\leq \|\beta_{\mathcal{A}^m}^m - \beta_{\mathcal{A}^m}^*\|_2 + \|\beta_{\mathcal{I}^m}^*\|_2 \\ &\leq \|(\mathbf{X}'_{\mathcal{A}^m} \mathbf{X}_{\mathcal{A}^m})^{-1} \mathbf{X}'_{\mathcal{A}^m} (\mathbf{X}\beta^* + \epsilon - \mathbf{X}_{\mathcal{A}^m} \beta_{\mathcal{A}^m}^*)\|_2 + \|\beta_{\mathcal{I}^m}^*\|_2 \\ &\leq \|(\mathbf{X}'_{\mathcal{A}^m} \mathbf{X}_{\mathcal{A}^m})^{-1} \mathbf{X}'_{\mathcal{A}^m} \mathbf{X}_{\mathcal{I}^m} \beta_{\mathcal{I}^m}^*\|_2 + \|(\mathbf{X}'_{\mathcal{A}^m} \mathbf{X}_{\mathcal{A}^m})^{-1} \mathbf{X}'_{\mathcal{A}^m} \epsilon\|_2 + \|\beta_{\mathcal{I}^m}^*\|_2 \\ &\leq (1 + \Delta + \frac{\theta_{s,s}}{c_-(s)}) \|\beta_{\mathcal{I}^m}^*\|_2 \\ &\leq \frac{1 + \Delta + \frac{\theta_{s,s}}{c_-(s)}}{(1-\Delta)n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)^2})} \delta_s^m |2n\mathcal{L}_n(\beta^0) - 2n\mathcal{L}_n(\beta^*)| \\ &\leq \frac{1 + \Delta + \frac{\theta_{s,s}}{c_-(s)}}{(1-\Delta)n(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)^2})} \delta_s^m \|\mathbf{y}\|_2^2. \end{aligned}$$

110 where the last second inequality follows from (9) and (10) .

111 Proof of Theorem 4.

112 **proof 3** For simplicity, denote $\mathcal{L}_{\mathcal{A}} = \min_{\beta_{\mathcal{I}}=0} \mathcal{L}_n(\beta)$, where $\mathcal{I} = (\mathcal{A})^c$. We need to bound $\log \mathcal{L}_{\mathcal{A}} - \log \mathcal{L}_{\mathcal{B}}$. It follows from

113 $1 - \frac{1}{x} \leq \log(x) \leq x - 1$ for $x > 0$, that

$$114 \quad \frac{\mathcal{L}_{\mathcal{A}} - \mathcal{L}_{\mathcal{B}}}{\mathcal{L}_{\mathcal{A}}} \leq \log \frac{\mathcal{L}_{\mathcal{A}}}{\mathcal{L}_{\mathcal{B}}} \leq \frac{\mathcal{L}_{\mathcal{A}} - \mathcal{L}_{\mathcal{B}}}{\mathcal{L}_{\mathcal{B}}}. \quad [13]$$

Let $(\hat{\beta}^s, \hat{\mathbf{d}}^s, \hat{\mathcal{A}}^s, \hat{\mathcal{I}}^s)$ be the output of Algorithm 1 with support size s . Firstly, consider the case when $s > s^*$. By Lemma 1, for any $s \geq s^*$,

$$\begin{aligned} \mathbb{P}(\hat{\mathcal{A}}^s \supseteq \mathcal{A}^*) &\geq 1 - \gamma(s, n, p, b^*) \\ &\geq 1 - O(\exp\{\log p - K_s \log p \log \log n\}) \\ &\geq 1 - O(p^{-\alpha}), \end{aligned}$$

for some constant $\alpha > 0$, and the last second inequality uses Condition (6) and Condition (7). Let $\hat{\mathcal{A}}^s = \mathcal{A}^* \cup \mathcal{B}^s$,

$$\begin{aligned} \mathcal{L}_{\mathcal{A}^*} - \mathcal{L}_{\hat{\mathcal{A}}^s} &= \frac{1}{2n} \mathbf{y}'(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{H}_{\hat{\mathcal{A}}^s} (\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{y} \\ &= \frac{1}{2n} \mathbf{y}'(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{X}_{\mathcal{B}^s} (\mathbf{X}'_{\mathcal{B}^s} (\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{X}_{\mathcal{B}^s})^{-1} \mathbf{X}'_{\mathcal{B}^s} (\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{y} \\ &= \frac{1}{2n} \epsilon' (\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{X}_{\mathcal{B}^s} (\mathbf{X}'_{\mathcal{B}^s} (\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \mathbf{X}_{\mathcal{B}^s})^{-1} \mathbf{X}'_{\mathcal{B}^s} (\mathbf{I} - \mathbf{H}_{\mathcal{A}^*}) \epsilon. \end{aligned}$$

Note that,

$$\begin{aligned}
& \mathbb{P} \left(\frac{1}{\sqrt{2n}} \|(\mathbf{X}'_{\mathcal{B}^s}(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*})\mathbf{X}_{\mathcal{B}^s})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{B}^s}(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*})\boldsymbol{\epsilon}\|_2 \geq t \right) \\
& \leq \sum_{j \in \mathcal{B}^s} \mathbb{P} \left(\frac{1}{\sqrt{2n}} \|e'_j(\mathbf{X}'_{\mathcal{B}^s}(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*})\mathbf{X}_{\mathcal{B}^s})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{B}^s}(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*})\boldsymbol{\epsilon}\|_2 \geq \frac{t}{\sqrt{|\mathcal{B}^s|}} \right) \\
& \leq 2p \exp \left\{ -\frac{2nt^2}{|\mathcal{B}^s|\sigma^2} \right\}.
\end{aligned}$$

Then with probability $1 - (2p)^{-\alpha}$, we have

$$\frac{1}{\sqrt{2n}} \|(\mathbf{X}'_{\mathcal{B}^s}(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*})\mathbf{X}_{\mathcal{B}^s})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{B}^s}(\mathbf{I} - \mathbf{H}_{\mathcal{A}^*})\boldsymbol{\epsilon}\|_2 \leq \sqrt{\frac{\sigma^2|\mathcal{B}^s|}{2n}(1 + \alpha) \log(2p)}. \quad [14]$$

Therefore, with probability $1 - O(p^{-\alpha})$,

$$\mathcal{L}_{\mathcal{A}^*} - \mathcal{L}_{\hat{\mathcal{A}}^s} \leq \frac{\sigma^2|\mathcal{B}^s|}{2n}(1 + \alpha) \log(2p). \quad [15]$$

115 Now turn to $\mathcal{L}_{\hat{\mathcal{A}}^s}$. Similar to (14), with probability $1 - O(p^{-\alpha})$

$$\begin{aligned}
\mathcal{L}_{\hat{\mathcal{A}}^s} &= \frac{1}{2n} \|\mathbf{y} - \mathbf{X}_{\hat{\mathcal{A}}^s} \hat{\boldsymbol{\beta}}_{\hat{\mathcal{A}}^s}\|_2^2 \\
&\geq \frac{1}{2n} \|\boldsymbol{\epsilon}\|_2^2 - \frac{1}{2n} \|\mathbf{X}_{\hat{\mathcal{A}}^s}(\hat{\boldsymbol{\beta}}_{\hat{\mathcal{A}}^s} - \boldsymbol{\beta}_{\hat{\mathcal{A}}^s}^*)\|_2^2 \\
&\geq \frac{\sigma^2}{2} \left(1 - \frac{1}{2n} \alpha \log(2p)\right) - \frac{\sigma^2 s}{2n} (1 + \alpha) \log(2p) \\
&> 0,
\end{aligned} \quad [16]$$

117 where the last second inequality uses the fact that $\|\boldsymbol{\epsilon}\|_2^2$ is sub-exponential random variables and the last inequality is induced by
118 Condition (7).

It follows from (13), (15) and (16) that with probability $1 - O(p^{-\alpha})$,

$$\log \frac{\mathcal{L}_{\mathcal{A}^*}}{\mathcal{L}_{\hat{\mathcal{A}}^s}} \leq \frac{\frac{\sigma^2|\mathcal{B}^s|}{n}(1 + \alpha) \log(2p)}{\frac{\sigma^2}{2} \left(1 - \frac{1}{2n} \alpha \log(2p)\right) - \frac{\sigma^2 s}{2n} (1 + \alpha) \log(2p)}.$$

Consequently,

$$\begin{aligned}
\text{SIC}(\mathcal{A}^*) - \text{SIC}(\hat{\mathcal{A}}^s) &= n \log \frac{\mathcal{L}_{\mathcal{A}^*}}{\mathcal{L}_{\hat{\mathcal{A}}^s}} - |\mathcal{B}^s| \log(p) \log \log n \\
&\leq O(|\mathcal{B}^s| \log(2p)) - |\mathcal{B}^s| \log(p) \log \log n \\
&< 0
\end{aligned}$$

119 for a sufficiently large n .

120 Next, consider the case when $s < s^*$. Denote $\mathcal{I}_1^s = \hat{\mathcal{I}}^s \cap \mathcal{A}^*$. Similar to the (2) in Lemma 1, with probability $1 - O(p^{-\alpha})$, for
121 some constants $0 < \Delta < \frac{1}{2}$, we have

$$\begin{aligned}
2\mathcal{L}_{\hat{\mathcal{A}}^s} &\geq (1 - \Delta) \left(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}\right) \|\boldsymbol{\beta}_{\mathcal{I}_1^s}^*\|_2^2 + \sigma^2, \\
2\mathcal{L}_{\hat{\mathcal{A}}^s} &\leq (1 + \Delta) c_+(s) \|\boldsymbol{\beta}_{\mathcal{I}_1^s}^*\|_2^2 + \sigma^2.
\end{aligned} \quad [17]$$

Denote $\mathbf{H}_{\mathcal{A}} = \mathbf{X}_{\mathcal{A}}(\mathbf{X}'_{\mathcal{A}}\mathbf{X}_{\mathcal{A}})^{-1}\mathbf{X}'_{\mathcal{A}}$. We have

$$2n(\mathcal{L}_{\hat{\mathcal{A}}^s} - \mathcal{L}_{\mathcal{A}^*}) \geq n(1 - \Delta) \left(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}\right) \|\boldsymbol{\beta}_{\mathcal{I}_1^s}^*\|_2^2 - \|\boldsymbol{\epsilon}^\top \mathbf{H}_{\mathcal{A}^*} \boldsymbol{\epsilon}\|_2^2.$$

We also have

$$\begin{aligned}
& \mathbb{P} \left(\|(\mathbf{X}'_{\mathcal{A}^*} \mathbf{X}_{\mathcal{A}^*})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{A}^*} \boldsymbol{\epsilon}\|_2 \geq t \right) \\
& \leq \sum_{j \in \mathcal{A}^*} \mathbb{P} \left(\|e'_j(\mathbf{X}'_{\mathcal{A}^*} \mathbf{X}_{\mathcal{A}^*})^{-\frac{1}{2}} \mathbf{X}'_{\mathcal{A}^*} \boldsymbol{\epsilon}\|_2 \geq \frac{t}{\sqrt{s^*}} \right) \\
& \leq 2p \exp \left\{ -\frac{t^2}{\sigma^2 s^*} \right\}.
\end{aligned}$$

123 Thus,

$$\begin{aligned}
& \mathbb{P} \left(\|\epsilon^\top \mathbf{H}_{\mathcal{A}^*} \epsilon\|_2^2 \geq n\Delta(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1^s}^*\|_2^2 \right) \\
& \leq 2 \exp \left\{ \log p - \frac{nK_{s,5}b^*}{s^*} \right\} \\
& \leq 2 \exp \{ \log p - K_{s,5} \log p \log \log n \} \\
& \leq O(p^{-\alpha}),
\end{aligned} \tag{18}$$

124 where $K_{s,5} = \frac{\Delta(c_-(s)^2 - \theta_{s,s}^2)}{c_-(s)\sigma^2}$ and the second inequality uses Condition (6).

Therefore, with probability $1 - O(p^{-\alpha})$, we have

$$2(\mathcal{L}_{\hat{\mathcal{A}}^s} - \mathcal{L}_{\mathcal{A}^*}) \geq (1 - 2\Delta)(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1^s}^*\|_2^2. \tag{19}$$

It follows from inequalities (13), (17) and (19) that with probability $1 - O(p^{-\alpha})$,

$$\log \frac{\mathcal{L}_{\hat{\mathcal{A}}^s}}{\mathcal{L}_{\mathcal{A}^*}} \geq \frac{(1 - 2\Delta)(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1^s}^*\|_2^2}{(1 + \Delta)c_+(s) \|\beta_{\mathcal{I}_1^s}^*\|_2^2 + \sigma^2}.$$

Consequently, for sufficiently large n ,

$$\begin{aligned}
\text{SIC}(\hat{\mathcal{A}}^s) - \text{SIC}(\mathcal{A}^*) &= n \log \frac{\mathcal{L}_{\hat{\mathcal{A}}^s}}{\mathcal{L}_{\mathcal{A}^*}} - (s^* - |\hat{\mathcal{A}}^s|) \log(p) \log \log n \\
&\geq nO(\min\{1, |\mathcal{I}_1^s|b^*\}) - |\mathcal{I}_1^s| \log(p) \log \log n \\
&> 0,
\end{aligned}$$

125 based on Condition (6) and Condition (7).

126 Therefore, information criterion $\text{SIC}(\hat{\mathcal{A}})$ attains the minimum at \mathcal{A}^* with probability $1 - O(p^{-\alpha})$.

127 Proof of Theorem 2.

128 **proof 4** Denote $\mathcal{L}_{\mathcal{A}} = \min_{\beta_{\mathcal{I}}=0} \mathcal{L}_n(\beta)$, where $\mathcal{I} = (\mathcal{A})^c$. First of all, consider $0 \leq s < s^*$. Since the loss function decreases by

129 at least τ_s at each iteration, Algorithm 1 stops before $O(\frac{\|\mathbf{y}\|_2^2}{\tau_s})$ iterations.

130 Next, consider $s^* \leq s \leq s_{\max}$, for a fix $s > 0$. Denote \mathcal{A}^m as the active set output by Algorithm 1 in m th iteration. Assuming $\mathcal{A}^m \not\supseteq \mathcal{A}^*$, by (9), we have

$$\mathcal{L}_{\mathcal{A}^m} - \mathcal{L}_{\mathcal{A}^{m+1}} \geq (1 - \delta_s)(\mathcal{L}_{\mathcal{A}^m} - \mathcal{L}_n(\beta^*)) \geq (1 - \delta_s)(1 - \Delta)(c_-(s) - \frac{\theta_{s,s}^2}{c_-(s)}) \|\beta_{\mathcal{I}_1^m}^*\|_2^2 > \tau_s,$$

131 so the difference between the m th loss function and the $(m+1)$ th loss function is bigger than the threshold τ_s . Thus, from

132 Theorem 3, we have $\mathcal{A}^m \supseteq \mathcal{A}^*$ after $O(\log \frac{\|\mathbf{y}\|_2^2}{s \log p \log \log n})$ iterations.

Now $\mathcal{A}^{m+1} \supseteq \mathcal{A}^*$, then $\mathcal{L}_{\mathcal{A}^{m+1}} = \epsilon' \mathbf{H}_{\mathcal{A}^{m+1}} \epsilon$ and $\mathcal{L}_{\mathcal{A}^m} = \epsilon' \mathbf{H}_{\mathcal{A}^m} \epsilon$. By (16), with probability $1 - O(p^{-\alpha})$, we have

$$\begin{aligned}
\mathcal{L}_{\mathcal{A}^m} &\leq \frac{\sigma^2}{2} + \frac{\sigma^2}{4n} \alpha \log(2p) + \frac{\sigma^2 s}{2n} (1 + \alpha) \log(2p), \\
\mathcal{L}_{\mathcal{A}^{m+1}} &\geq \frac{\sigma^2}{2} - \frac{\sigma^2}{4n} \alpha \log(2p) - \frac{\sigma^2 s}{2n} (1 + \alpha) \log(2p).
\end{aligned}$$

Thus, for sufficient large n , we have

$$\mathcal{L}_{\mathcal{A}^m} - \mathcal{L}_{\mathcal{A}^{m+1}} \leq \frac{\sigma^2}{2n} \alpha \log(2p) + \frac{\sigma^2 s}{n} (1 + \alpha) \log(2p) \leq \tau_s.$$

133 Therefore, for $s^* \leq s \leq s_{\max}$, Algorithm 1 terminates after $O(\log \frac{\|\mathbf{y}\|_2^2}{s \log p \log \log n})$ iterations.

Now we analyze the computational complexity of Algorithm 2 for a given active set size s . Computing ξ and ζ takes $O(np + ns)$ steps. Finding the smallest (largest) s values in ξ (ζ) takes $O(p)$ steps via Hoare's selection algorithm (12). Because the procedure repeats k_{\max} , $O(k_{\max}p)$ steps are demanded. In Algorithm 2, the splicing method iterates at most s times. Therefore the computational complexity of Algorithm 1 is

$$O \left(\log \frac{\|\mathbf{y}\|_2^2}{s \log p \log \log n} \mathbb{I}(s^* \leq s) + \frac{\|\mathbf{y}\|_2^2}{\tau_s} \mathbb{I}(s^* > s) \right) \cdot O(nsp + ns^2 + k_{\max}sp),$$

where $\mathbb{I}(\cdot)$ is a indicator function. Since s varies from 1 to s_{\max} in Algorithm 3, the total computational complexity is

$$\begin{aligned}
& O \left(\log \frac{\|\mathbf{y}\|_2^2}{\log p \log \log n} (nps_{\max}^2 + ns_{\max}^3 + k_{\max}ps_{\max}^2) + \frac{n\|\mathbf{y}\|_2^2}{\log p \log \log n} (nps^* + n(s^*)^2 + k_{\max}ps^*) \right) \\
& \leq O \left((s_{\max} \log \frac{\|\mathbf{y}\|_2^2}{\log p \log \log n} + \frac{n\|\mathbf{y}\|_2^2}{\log p \log \log n}) (nps_{\max} + ns_{\max}^2 + k_{\max}ps_{\max}) \right).
\end{aligned}$$

134 **References**

- 135 1. H Hazimeh, R Mazumder, Fast best subset selection: Coordinate descent and local combinatorial optimization algorithms.
136 *Oper. Res.* (2020).
- 137 2. C Wen, A Zhang, S Quan, X Wang, Bess: An r package for best subset selection in linear, logistic and cox proportional
138 hazards models. *J. Stat. Softw.* **94**, 1–24 (2020).
- 139 3. J Friedman, T Hastie, R Tibshirani, Regularization paths for generalized linear models via coordinate descent. *J. Stat.*
140 *Software, Articles* **33**, 1–22 (2010).
- 141 4. P Breheny, J Huang, Coordinate descent algorithms for nonconvex penalized regression, with applications to biological
142 feature selection. *Ann. Appl. Stat.* **5**, 232–253 (2011).
- 143 5. H Hazimeh, R Mazumder, *L0Learn: Fast Algorithms for Best Subset Selection*, (2019) R package version 1.2.0.
- 144 6. C Wen, A Zhang, S Quan, X Wang, *BeSS: Best Subset Selection for Sparse Generalized Linear Model and Cox Model*,
145 (2017) R package version 1.0.2.
- 146 7. J Huang, Y Jiao, Y Liu, X Lu, A constructive approach to l_0 penalized regression. *J. Mach. Learn. Res.* **19**, 1–37 (2018).
- 147 8. L Wang, Y Kim, R Li, Calibrating non-convex penalized regression in ultra-high dimension. *Annals statistics* **41**, 2505
148 (2013).
- 149 9. DF Saldana, Y Feng, Sis: an r package for sure independence screening in ultrahigh dimensional statistical models. *J.*
150 *Stat. Softw.* **83**, 1–25 (2018).
- 151 10. H Zou, T Hastie, R Tibshirani, , et al., On the “degrees of freedom” of the lasso. *The Annals Stat.* **35**, 2173–2192 (2007).
- 152 11. J Kiefer, Sequential minimax search for a maximum. *Proc. Am. mathematical society* **4**, 502–506 (1953).
- 153 12. CAR Hoare, Algorithm 65: Find. *Commun. ACM* **4**, 321–322 (1961).