

Supplemental material to:
 “The matching pursuit algorithm revisited:
 A variant for big data and new stopping rules”

Fangyao Li^a, Christopher M. Triggs^a, Bogdan Dumitrescu^b, and Ciprian Doru
 Giurcăneanu^a

^aDept. of Statistics, Univ. of Auckland, New Zealand

^bUniv. Politehnica of Bucharest, Romania

Contents

1	Additional information on the simulation study	2
1.1	MISE tables	2
1.1.1	Experiments for $n = 20$	2
1.1.2	Experiments for $n = 100$	4
1.1.3	Experiments for $n = 10000$	6
1.2	Median number of iterations	8
1.2.1	Experiments for $n = 20$	8
1.2.2	Experiments for $n = 100$	10
1.2.3	Experiments for $n = 10000$	12
1.3	Median number of predictors	14
1.3.1	Experiments for $n = 20$	14
1.3.2	Experiments for $n = 100$	16
1.3.3	Experiments for $n = 10000$	18
1.4	Score plots	20
2	Additional information on experiments with air pollution data	23
2.1	Location of the four sites where the concentrations of the air pollutants are measured	23
2.2	Detailed information on the experiments with air pollution data collected from Patumahoe site	24
2.2.1	The number of iterations and the number of predictors	24
2.2.2	Statistics on how many times each predictor is selected in 100 runs	27
2.2.3	How many times each predictor from the constrained set is selected in 100 runs: The case when the total number of predictors is $p_n = 68$ (ConSet) versus the case when the total number of predictors is $p_n = 1464$ (FullSet)	28
2.2.4	Statistics for EgMDL ₁ [*] on how many times each predictor from the FullSet ($p_n = 1464$) not included in the ConSet ($p_n = 68$) is selected in 100 runs	30

1 Additional information on the simulation study

In this section, we provide the experimental results obtained with simulated data. The procedure for generating the data is presented in [1, Sec. 5.1].

1.1 MISE tables

We use the mean integrated squared error (MISE) to evaluate the performance of the stopping rules. The formula of MISE is given in [1, Eq. (27)]. For each sample size $n \in \{20, 100, 10000\}$, the results are shown in four different tables; each table corresponds to one of the four data models defined in [1, Sec. 5.1]. Note that the size of the dictionaries is the same in all the experiments ($p_n = 100$). For each row of the table, the minimum MISE is in red, and the results which are within a range of 5% from the minimum value on that row are shown in bold.

1.1.1 Experiments for $n = 20$

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ ⁺	ESC ₂ ⁺	gMDL ₁	gMDL ₂	E _g MDL ₁	E _g MDL ₂	E _g MDL ₁ ⁺	E _g MDL ₂ ⁺	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.1288	0.1294	0.1286	0.1294	0.1294	0.1294	0.1367	0.1291	0.1291	0.1289	0.1290	0.1449	0.1560	0.1290	0.1291	0.1289	0.1289	0.1363	0.1340	0.1291	0.1294	0.1294	0.1358
(0.00, 0.20)	1.1500	1.1682	1.1263	1.1682	1.1682	1.1682	1.0857	1.1674	1.1674	1.1663	1.1663	0.3811	0.3714	1.1670	1.1670	1.1662	1.1663	0.4110	0.4337	1.1670	1.1682	1.1682	0.6973
(0.75, 8.00)	0.1188	0.1306	0.1115	0.1306	0.1306	0.1306	0.1230	0.1245	0.1248	0.1182	0.1183	0.0818	0.0818	0.1241	0.1248	0.1178	0.1184	0.0812	0.0811	0.1250	0.1306	0.1306	0.0932
(0.75, 0.20)	3.0117	3.4640	2.5065	3.4640	3.4640	3.4640	2.3669	3.4531	3.4574	3.4503	3.4508	0.8154	0.8054	3.4517	3.4531	3.4491	3.4505	0.8325	0.8125	3.4531	3.4627	3.4627	1.7325
Case 2																							
(0.00, 8.00)	0.3451	0.3639	0.3360	0.3639	0.3639	0.3639	0.3515	0.3447	0.3447	0.3308	0.3309	0.3350	0.3348	0.3444	0.3445	0.3310	0.3315	0.3328	0.3324	0.3473	0.3639	0.3639	0.2948
(0.00, 0.20)	2.7134	3.1700	2.4078	3.1700	3.1700	3.1700	2.4100	3.1295	3.1317	2.9593	2.9607	1.2485	1.2402	3.1297	3.1312	2.9583	2.9595	1.2699	1.2792	3.1312	3.1699	3.1699	1.3109
(0.75, 8.00)	0.3888	0.4470	0.3577	0.4470	0.4470	0.4470	0.3565	0.3613	0.3616	0.3519	0.3520	0.2808	0.2810	0.3617	0.3632	0.3553	0.3558	0.2816	0.2814	0.3632	0.4407	0.4408	0.1733
(0.75, 0.20)	6.9660	9.5772	6.0873	9.5772	9.5772	9.5772	5.0611	8.9919	8.9949	8.7405	8.7483	3.1070	3.1031	8.9986	9.0011	8.7428	8.8301	3.1696	3.2154	9.0023	9.4961	9.4961	3.2113
Case 3																							
(0.00, 8.00)	0.1971	0.1987	0.1967	0.1987	0.1987	0.1987	0.2078	0.1987	0.1974	0.1980	0.1979	0.2215	0.2213	0.1991	0.1990	0.1978	0.1978	0.2188	0.2178	0.1974	0.1987	0.1987	0.2047
(0.00, 0.20)	1.4421	1.5219	1.3653	1.5219	1.5219	1.5219	1.2699	1.5049	1.5050	1.4989	1.5006	0.7628	0.7553	1.5002	1.5048	1.4990	1.5006	0.7888	0.7952	1.5048	1.5219	1.5219	0.9504
(0.75, 8.00)	0.1811	0.2277	0.1633	0.2277	0.2277	0.2277	0.1841	0.1998	0.2001	0.1741	0.1768	0.1199	0.1200	0.2003	0.2003	0.1753	0.1813	0.1205	0.1206	0.2003	0.2274	0.2274	0.1375
(0.75, 0.20)	4.1145	5.8723	3.4212	5.8723	5.8723	5.8723	3.3030	5.7323	5.7736	5.6917	5.7019	1.3266	1.2770	5.7370	5.7751	5.6925	5.6953	1.3824	1.3991	5.7756	5.8717	5.8717	3.1255

Table 1: MISE: Model 1 (low-dimensional)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ ⁺	ESC ₂ ⁺	gMDL ₁	gMDL ₂	E _g MDL ₁	E _g MDL ₂	E _g MDL ₁ ⁺	E _g MDL ₂ ⁺	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.0113	0.0113	0.0112	0.0113	0.0113	0.0113	0.0111	0.0113	0.0113	0.0113	0.0113	0.0103	0.0102	0.0113	0.0113	0.0113	0.0113	0.0104	0.0104	0.0113	0.0113	0.0113	0.0108
(0.00, 0.20)	0.0363	0.0373	0.0357	0.0373	0.0373	0.0373	0.0332	0.0373	0.0373	0.0372	0.0372	0.0121	0.0120	0.0373	0.0373	0.0372	0.0372	0.0129	0.0136	0.0373	0.0373	0.0373	0.0235
(0.75, 8.00)	0.0357	0.0354	0.0361	0.0354	0.0354	0.0354	0.0483	0.0356	0.0355	0.0357	0.0357	0.0522	0.0537	0.0356	0.0355	0.0357	0.0357	0.0498	0.0493	0.0355	0.0354	0.0354	0.0395
(0.75, 0.20)	0.2400	0.2807	0.2154	0.2807	0.2807	0.2807	0.1782	0.2802	0.2802	0.2799	0.2799	0.0720	0.0707	0.2802	0.2802	0.2798	0.2799	0.0750	0.0754	0.0355	0.2807	0.2807	0.1545
Case 2																							
(0.00, 8.00)	0.0168	0.0175	0.0162	0.0175	0.0175	0.0175	0.0159	0.0174	0.0174	0.0172	0.0172	0.0145	0.0145	0.0174	0.0174	0.0172	0.0172	0.0145	0.0146	0.0173	0.0175	0.0175	0.0140
(0.00, 0.20)	0.0871	0.0977	0.0766	0.0977	0.0977	0.0977	0.0778	0.0968	0.0968	0.0931	0.0934	0.0378	0.0374	0.0968	0.0968	0.0929	0.0933	0.0388	0.0389	0.0968	0.0977	0.0977	0.0539
(0.75, 8.00)	0.0703	0.0736	0.0706	0.0736	0.0736	0.0736	0.0722	0.0680	0.0680	0.0678	0.0678	0.0738	0.0741	0.0689	0.0689	0.0680	0.0680	0.0736	0.0736	0.0690	0.0690	0.0690	0.0660
(0.75, 0.20)	0.7123	1.0147	0.5931	1.0147	1.0147	1.0147	0.4068	0.9561	0.9637	0.9001	0.9037	0.2432	0.2421	0.9569	0.9576	0.9000	0.9036	0.2497	0.2524	0.9576	1.0147	1.0147	0.3768
Case 3																							
(0.00, 8.00)	0.0123	0.0124	0.0122	0.0124	0.0124	0.0124	0.0119	0.0124	0.0124	0.0124	0.0124	0.0114	0.0114	0.0124	0.0124	0.0124	0.0124	0.0114	0.0114	0.0124	0.0124	0.0124	0.0108
(0.00, 0.20)	0.0466	0.0495	0.0440	0.0495	0.0495	0.0495	0.0390	0.0492	0.0492	0.0489	0.0489	0.0218	0.0215	0.0492	0.0492	0.0489	0.0490	0.0228	0.0230	0.0492	0.0495	0.0495	0.0242
(0.75, 8.00)	0.0406	0.0400	0.0409	0.0400	0.0400	0.0400	0.0514	0.0410	0.0410	0.0410	0.0412	0.0548	0.0547	0.0410	0.0410	0.0414	0.0413	0.0533	0.0541	0.0409	0.0400	0.0400	0.0474
(0.75, 0.20)	0.3290	0.4178	0.2873	0.4178	0.4178	0.4178	0.2111	0.4128	0.4138	0.4087	0.4091	0.1243	0.1216	0.4115	0.4130	0.4078	0.4090	0.1280	0.1284	0.4132	0.4177	0.4177	0.1738

Table 2: MISE: Model 2 (high-dimensional, small equal coefficients)

(ω, ς^2)	AIC_C	KIC	KIC_C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.86	0.86	0.85	0.86	0.86	0.86	0.92	0.86	0.86	0.86	0.86	0.89	0.94	0.86	0.86	0.86	0.86	0.84	0.86	0.86	0.86	0.86	0.87
(0.00, 0.20)	5.75	5.85	5.62	5.85	5.85	5.85	5.40	5.84	5.84	5.84	5.84	1.99	1.94	5.84	5.84	5.84	5.84	2.22	2.33	5.84	5.85	5.85	3.39
(0.75, 8.00)	1.24	1.29	1.23	1.29	1.29	1.29	1.70	1.27	1.27	1.25	1.26	1.44	1.44	1.27	1.27	1.25	1.26	1.42	1.42	1.27	1.29	1.29	1.27
(0.75, 0.20)	19.81	22.27	17.96	22.27	22.27	22.27	15.27	22.25	22.25	22.16	22.16	5.99	5.80	22.22	22.25	22.14	22.16	6.14	6.14	22.25	22.27	22.27	10.72
Case 2																							
(0.00, 8.00)	1.57	1.65	1.53	1.65	1.65	1.65	1.67	1.54	1.55	1.52	1.52	1.48	1.48	1.55	1.55	1.52	1.52	1.48	1.48	1.55	1.65	1.65	1.37
(0.00, 0.20)	12.42	13.61	11.57	13.61	13.61	13.61	10.37	12.96	12.96	12.75	12.77	6.69	6.64	12.91	12.91	12.71	12.73	6.75	6.83	12.92	13.61	13.61	5.71
(0.75, 8.00)	4.42	4.93	4.01	4.93	4.93	4.93	4.16	4.06	4.07	3.94	3.94	3.68	3.68	4.08	4.19	3.95	3.95	3.69	3.70	4.19	4.91	4.91	2.58
(0.75, 0.20)	57.30	79.92	47.86	79.92	79.92	79.92	37.59	72.11	77.73	66.44	66.53	22.31	22.27	71.58	72.45	66.50	66.48	22.93	23.12	76.72	79.91	79.91	26.99
Case 3																							
(0.00, 8.00)	1.09	1.11	1.08	1.11	1.11	1.11	1.20	1.11	1.11	1.10	1.09	1.15	1.16	1.11	1.11	1.10	1.09	1.12	1.12	1.11	1.11	1.11	1.11
(0.00, 0.20)	7.01	7.53	6.59	7.53	7.53	7.53	5.30	7.49	7.49	7.45	7.47	3.40	3.34	7.49	7.49	7.45	7.46	3.53	3.57	7.49	7.53	7.53	3.85
(0.75, 8.00)	1.65	1.76	1.58	1.76	1.76	1.76	2.01	1.63	1.63	1.59	1.60	1.69	1.68	1.63	1.62	1.60	1.60	1.65	1.65	1.62	1.76	1.76	1.78
(0.75, 0.20)	26.06	32.32	21.87	32.32	32.32	32.32	19.88	31.72	31.76	31.07	31.17	9.17	9.06	31.65	31.75	30.93	31.16	9.83	9.78	31.80	32.31	32.31	16.67

Table 3: MISE: Model 3 (high-dimensional, decaying coefficients)

(ω, ς^2)	AIC_C	KIC	KIC_C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	5.56	5.57	5.55	5.57	5.57	5.57	5.51	5.57	5.57	5.57	5.57	5.14	5.15	5.57	5.57	5.57	5.57	5.19	5.18	5.57	5.57	5.57	5.33
(0.00, 0.20)	18.06	18.41	17.66	18.41	18.41	18.41	16.52	18.40	18.40	18.39	18.39	6.73	6.65	18.39	18.40	18.36	18.39	7.12	7.35	18.40	18.41	18.41	9.73
(0.75, 8.00)	12.54	12.52	12.75	12.52	12.52	12.52	17.18	12.48	12.47	12.56	12.53	18.37	18.82	12.50	12.47	12.57	12.53	17.02	17.10	12.47	12.52	12.52	13.79
(0.75, 0.20)	111.24	131.16	96.41	131.16	131.16	131.16	72.15	130.68	130.87	129.98	130.08	30.58	30.34	130.67	130.77	129.97	130.07	32.64	32.94	130.80	131.15	131.15	62.70
Case 2																							
(0.00, 8.00)	8.20	8.47	7.90	8.47	8.47	8.47	7.81	8.24	8.24	8.11	8.10	6.83	6.84	8.24	8.25	8.11	8.14	6.83	6.84	8.25	8.47	8.47	6.58
(0.00, 0.20)	32.98	38.01	30.22	38.01	38.01	38.01	30.24	35.85	36.90	34.09	34.21	18.08	17.90	35.86	37.08	34.27	34.30	18.28	18.31	37.08	38.02	38.02	16.18
(0.75, 8.00)	24.99	26.97	24.24	26.97	26.97	26.97	26.09	25.22	25.30	23.63	23.69	26.29	26.38	25.22	25.20	23.60	23.65	25.81	25.80	25.20	26.91	26.91	22.58
(0.75, 0.20)	296.10	367.98	239.26	367.98	367.98	367.98	222.31	351.78	352.58	346.67	346.73	116.94	116.57	352.52	352.55	346.77	346.78	118.62	119.09	352.62	367.87	367.87	149.52
Case 3																							
(0.00, 8.00)	6.07	6.17	6.00	6.17	6.17	6.17	5.79	6.14	6.14	6.14	6.14	5.48	5.49	6.14	6.14	6.14	6.14	5.49	5.48	6.14	6.17	6.17	5.42
(0.00, 0.20)	22.78	24.84	21.52	24.84	24.84	24.84	19.87	24.73	24.73	24.73	24.73	12.18	11.96	24.72	24.73	24.73	24.73	12.57	12.69	24.73	24.84	24.84	13.58
(0.75, 8.00)	14.44	14.51	14.54	14.51	14.51	14.51	18.07	14.58	14.58	14.62	14.68	19.68	19.87	14.60	14.57	14.60	14.66	19.39	19.25	14.55	14.51	14.51	17.01
(0.75, 0.20)	124.94	166.72	105.57	166.72	166.72	166.72	88.89	163.27	163.57	160.06	160.33	47.98	46.29	163.47	163.61	160.00	160.12	49.08	49.79	163.61	166.71	166.71	78.59

Table 4: MISE: Model 4 (high-dimensional, slowly decaying coefficients)

1.1.2 Experiments for $n = 100$

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.0198	0.1197	0.0136	0.0247	0.0326	0.0097	0.0104	0.0099	0.0098	0.0092	0.0091	0.0088	0.0088	0.0101	0.0100	0.0093	0.0092	0.0088	0.0088	0.0099	0.0130	0.0127	0.0129
(0.00, 0.20)	0.7000	4.6855	0.4470	0.8136	0.9396	0.2569	0.2966	0.9376	0.9985	1.1047	1.3267	0.2562	0.2687	0.9450	1.0289	1.1185	1.3472	0.2522	0.2460	1.1692	5.1333	5.6680	0.4250
(0.75, 8.00)	0.0369	0.1513	0.0210	0.0200	0.0180	0.0126	0.0145	0.0134	0.0133	0.0129	0.0129	0.0130	0.0131	0.0137	0.0137	0.0130	0.0130	0.0129	0.0129	0.0136	0.0143	0.0141	0.0209
(0.75, 0.20)	1.2739	5.7434	0.6671	0.6084	0.5120	0.3089	0.4498	3.0562	3.0106	4.1354	4.2554	0.3290	0.3442	3.0970	3.0709	4.2935	4.5241	0.3101	0.3026	3.5629	8.3796	8.4950	0.5183
Case 2																							
(0.00, 8.00)	0.9684	1.2011	0.9013	0.9735	0.9987	0.7588	0.4887	0.7540	0.7541	0.7483	0.7506	0.6117	0.6119	0.7556	0.7557	0.7524	0.7526	0.6130	0.6132	0.7559	0.8135	0.8134	0.2436
(0.00, 0.20)	8.4636	14.6361	6.7935	7.9764	8.2719	3.8526	2.5644	6.2392	6.2362	6.1821	6.2122	3.4752	3.5105	6.2745	6.2694	6.2545	6.2755	3.5022	3.5337	6.2845	8.2143	8.2059	1.2963
(0.75, 8.00)	2.9964	3.1696	2.8282	2.8880	2.8950	2.2789	1.8946	2.4209	2.4213	2.3820	2.3849	2.0724	2.0730	2.4262	2.4272	2.3945	2.3955	2.0755	2.0756	2.4283	2.5094	2.5098	0.1291
(0.75, 0.20)	27.9573	33.7314	21.4106	25.6019	25.8813	11.5545	8.0098	20.3783	20.2281	21.0777	20.9260	11.0387	11.0553	20.5232	20.6722	21.2022	21.2490	11.1062	11.1112	21.0016	27.9863	27.9293	3.3457
Case 3																							
(0.00, 8.00)	0.1832	0.3496	0.1576	0.1739	0.1783	0.1079	0.0803	0.1243	0.1244	0.1187	0.1180	0.0980	0.0981	0.1249	0.1249	0.1199	0.1198	0.0983	0.0984	0.1248	0.1393	0.1394	0.0870
(0.00, 0.20)	2.7677	8.3255	2.0594	2.4475	2.6466	1.0509	0.6973	2.3019	2.3024	2.3299	2.3430	1.0505	1.0383	2.3192	2.3417	2.3483	2.3773	1.0656	1.0528	2.3461	4.1942	4.2942	1.6430
(0.75, 8.00)	0.4833	0.6329	0.4217	0.4612	0.4631	0.2961	0.2457	0.3485	0.3486	0.3310	0.3313	0.2782	0.2778	0.3510	0.3522	0.3376	0.3378	0.2805	0.2806	0.3521	0.3808	0.3807	0.0879
(0.75, 0.20)	7.2871	13.1273	5.1037	6.3559	5.8106	2.1596	1.5896	6.5055	6.4625	6.6830	6.5782	2.1465	2.1461	6.6102	6.5466	6.9302	6.9251	2.1870	2.1642	6.6084	11.6053	11.7188	2.2397

Table 5: MISE: Model 1 (low-dimensional)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.0061	0.0063	0.0066	0.0062	0.0062	0.0073	0.0100	0.0061	0.0060	0.0060	0.0059	0.0088	0.0093	0.0061	0.0060	0.0060	0.0059	0.0086	0.0091	0.0060	0.0055	0.0055	0.0062
(0.00, 0.20)	0.0239	0.1497	0.0181	0.0248	0.0310	0.0120	0.0100	0.0309	0.0343	0.0347	0.0410	0.0101	0.0101	0.0313	0.0346	0.0356	0.0414	0.0103	0.0102	0.0365	0.1528	0.1738	0.0152
(0.75, 8.00)	0.0083	0.0155	0.0087	0.0094	0.0096	0.0106	0.0420	0.0093	0.0091	0.0093	0.0091	0.0118	0.0116	0.0093	0.0091	0.0092	0.0091	0.0118	0.0115	0.0091	0.0087	0.0086	0.0096
(0.75, 0.20)	0.1167	0.4954	0.0744	0.0622	0.0593	0.0568	0.0650	0.2324	0.2095	0.3152	0.3098	0.0558	0.0595	0.2394	0.2169	0.3235	0.3173	0.0555	0.0554	0.2300	0.6845	0.7078	0.0715
Case 2																							
(0.00, 8.00)	0.0317	0.0352	0.0290	0.0344	0.0347	0.0311	0.0204	0.0267	0.0267	0.0272	0.0272	0.0235	0.0237	0.0268	0.0268	0.0272	0.0273	0.0237	0.0238	0.0268	0.0296	0.0297	0.0128
(0.00, 0.20)	0.2558	0.4373	0.1915	0.2403	0.2619	0.1174	0.0806	0.1780	0.1783	0.1782	0.1784	0.1104	0.1108	0.1791	0.1794	0.1801	0.1792	0.1114	0.1119	0.1795	0.2451	0.2445	0.0475
(0.75, 8.00)	0.1933	0.2065	0.1763	0.1914	0.1941	0.1598	0.1054	0.1481	0.1481	0.1495	0.1496	0.1253	0.1257	0.1481	0.1482	0.1497	0.1497	0.1258	0.1257	0.1482	0.1605	0.1605	0.0659
(0.75, 0.20)	2.2498	2.6826	1.7057	2.0762	2.1435	0.9421	0.6067	1.6292	1.6156	1.7112	1.6886	0.8774	0.8787	1.6898	1.6677	1.7229	1.7155	0.8917	0.8943	1.6755	2.2373	2.2386	0.2929
Case 3																							
(0.00, 8.00)	0.0110	0.0134	0.0111	0.0126	0.0128	0.0125	0.0111	0.0111	0.0111	0.0111	0.0111	0.0110	0.0111	0.0111	0.0111	0.0111	0.0111	0.0111	0.0111	0.0111	0.0112	0.0112	0.0106
(0.00, 0.20)	0.0801	0.2248	0.0621	0.0752	0.0780	0.0316	0.0234	0.0702	0.0702	0.0720	0.0725	0.0312	0.0311	0.0708	0.0708	0.0743	0.0734	0.0318	0.0315	0.0719	0.1270	0.1285	0.0464
(0.75, 8.00)	0.0402	0.0558	0.0349	0.0410	0.0444	0.0313	0.0417	0.0319	0.0319	0.0320	0.0320	0.0290	0.0291	0.0319	0.0320	0.0321	0.0320	0.0291	0.0291	0.0320	0.0335	0.0335	0.0436
(0.75, 0.20)	0.7044	1.2702	0.5045	0.5699	0.6034	0.2036	0.1502	0.6260	0.6201	0.6621	0.6453	0.2038	0.2025	0.6333	0.6227	0.6650	0.6559	0.2060	0.2053	0.6344	1.1067	1.1161	0.4076

Table 6: MISE: Model 2 (high-dimensional, small equal coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV	
Case 1																								
(0.00, 8.00)	0.20	0.67	0.18	0.27	0.30	0.24	0.37	0.18	0.18	0.18	0.19	0.24	0.24	0.18	0.18	0.18	0.18	0.23	0.24	0.18	0.19	0.19	0.22	
(0.00, 0.20)	3.17	20.78	2.21	2.13	2.19	1.21	1.38	4.70	4.56	5.28	6.28	1.25	1.28	4.76	4.62	5.33	6.38	1.23	1.19	5.37	23.62	27.26	1.70	
(0.75, 8.00)	0.34	1.05	0.29	0.32	0.36	0.34	0.60	0.28	0.28	0.29	0.28	0.37	0.37	0.28	0.28	0.28	0.28	0.37	0.37	0.28	0.28	0.28	0.28	0.37
(0.75, 0.20)	9.75	41.01	5.29	5.86	5.94	2.94	3.99	19.90	19.52	27.22	27.53	3.03	3.22	20.07	20.39	27.52	28.92	2.83	2.85	22.22	54.64	55.69	5.65	
Case 2																								
(0.00, 8.00)	4.32	5.36	3.95	4.70	4.90	3.63	2.36	3.44	3.45	3.45	3.45	2.83	2.83	3.45	3.45	3.45	3.45	2.83	2.84	3.45	3.64	3.64	3.64	1.11
(0.00, 0.20)	38.66	72.28	31.19	36.79	38.89	19.34	13.81	27.85	27.94	27.89	27.82	17.47	17.49	28.12	28.12	28.23	28.06	17.62	17.66	28.10	38.06	38.16	38.16	6.95
(0.75, 8.00)	16.55	17.59	15.46	16.07	16.47	12.81	8.67	13.48	13.48	13.30	13.30	11.07	11.09	13.55	13.56	13.29	13.39	11.09	11.10	13.56	14.18	14.19	14.19	2.24
(0.75, 0.20)	165.99	205.77	132.95	151.88	156.85	69.93	39.59	138.62	134.35	139.55	139.35	61.98	62.12	139.44	135.14	140.15	141.20	62.62	62.89	139.94	175.04	173.75	173.75	16.35
Case 3																								
(0.00, 8.00)	1.07	1.88	0.94	1.26	1.36	0.83	0.66	0.85	0.85	0.84	0.84	0.72	0.72	0.85	0.85	0.84	0.85	0.72	0.72	0.85	0.92	0.92	0.92	0.70
(0.00, 0.20)	12.86	39.34	9.73	11.18	14.83	4.82	3.45	11.12	11.08	12.25	11.59	4.86	4.82	11.31	11.26	12.09	11.66	4.90	4.85	11.33	21.95	22.74	22.74	7.24
(0.75, 8.00)	3.76	4.90	3.45	3.80	3.86	2.51	2.02	2.89	2.90	2.84	2.85	2.34	2.35	2.91	2.91	2.85	2.86	2.35	2.35	2.92	3.15	3.15	3.15	1.80
(0.75, 0.20)	50.14	95.06	39.07	43.13	44.00	18.09	11.87	46.06	44.87	49.61	49.69	17.72	17.64	46.53	45.24	49.89	50.56	18.16	18.04	46.44	78.18	77.88	77.88	21.01

Table 7: MISE: Model 3 (high-dimensional, decaying coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV	
Case 1																								
(0.00, 8.00)	1.97	2.83	2.11	2.42	2.58	2.97	4.40	2.03	2.01	2.01	1.99	2.95	3.12	2.02	2.00	2.00	1.98	2.89	3.05	2.00	1.92	1.92	1.92	2.25
(0.00, 0.20)	12.00	78.24	8.82	9.19	13.96	5.08	5.09	16.38	17.27	18.25	23.10	5.14	5.12	16.65	17.69	18.56	23.60	5.17	5.11	19.81	81.80	88.91	88.91	7.45
(0.75, 8.00)	2.94	6.06	2.96	3.02	3.23	3.72	11.91	3.13	3.09	3.14	3.10	4.22	4.17	3.12	3.09	3.13	3.09	4.20	4.13	3.09	2.97	2.95	2.95	3.32
(0.75, 0.20)	48.13	194.83	27.42	21.56	19.97	23.00	26.96	89.33	83.05	132.40	118.84	22.79	24.32	90.28	85.63	134.97	137.82	22.36	22.23	97.41	295.25	299.78	299.78	28.92
Case 2																								
(0.00, 8.00)	15.78	18.26	14.34	17.62	17.81	15.69	10.99	13.07	13.08	13.23	13.28	11.82	11.83	13.10	13.13	13.26	13.28	11.83	11.83	13.22	14.17	14.19	14.19	6.34
(0.00, 0.20)	145.71	255.05	116.35	142.00	145.97	65.09	45.61	104.78	104.72	104.71	104.51	60.53	60.65	105.58	105.88	108.45	105.30	60.83	60.79	105.51	157.79	153.67	153.67	16.50
(0.75, 8.00)	88.41	95.35	80.69	86.76	88.38	70.36	43.90	69.61	69.65	69.88	69.91	57.80	57.84	69.89	69.91	70.06	70.06	58.12	58.16	69.91	73.94	73.95	73.95	25.05
(0.75, 0.20)	1080.09	1348.46	877.10	1008.66	1048.76	430.86	260.99	826.45	827.38	859.00	857.75	360.27	361.87	838.42	837.09	871.15	867.14	366.83	366.01	843.73	1129.92	1129.47	1129.47	108.30
Case 3																								
(0.00, 8.00)	4.67	6.97	4.43	6.23	6.45	6.01	5.08	4.30	4.31	4.33	4.36	4.57	4.57	4.30	4.31	4.36	4.37	4.53	4.57	4.31	4.60	4.62	4.62	4.67
(0.00, 0.20)	50.50	145.14	35.75	44.02	54.14	18.46	13.20	40.02	40.18	40.45	40.81	18.44	18.33	40.52	40.63	42.05	41.27	18.66	18.57	40.57	82.00	82.89	82.89	28.42
(0.75, 8.00)	17.79	24.49	15.67	17.29	17.93	12.98	12.98	13.50	13.57	13.46	13.50	11.94	11.98	13.57	13.60	13.53	13.55	11.94	11.99	13.60	14.60	14.61	14.61	15.48
(0.75, 0.20)	250.85	503.84	178.78	197.10	187.31	82.91	59.98	215.95	212.79	229.78	224.71	82.56	82.38	219.75	214.44	238.63	237.51	83.49	83.52	219.52	433.51	433.78	433.78	121.93

Table 8: MISE: Model 4 (high-dimensional, slowly decaying coefficients)

1.1.3 Experiments for $n = 10000$

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV	
Case 1																								
(0.00, 8.00)	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0008
(0.00, 0.20)	0.0046	0.0034	0.0034	0.0024	0.0025	0.0028	0.0033	0.0025	0.0025	0.0025	0.0024	0.0026	0.0026	0.0025	0.0025	0.0025	0.0025	0.0026	0.0026	0.0024	0.0026	0.0025	0.0696	
(0.75, 8.00)	0.0003	0.0002	0.0002	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001	0.0104	
(0.75, 0.20)	0.0101	0.0057	0.0057	0.0032	0.0033	0.0035	0.0040	0.0034	0.0033	0.0033	0.0034	0.0034	0.0034	0.0034	0.0034	0.0034	0.0034	0.0034	0.0034	0.0032	0.0034	0.0034	0.1298	
Case 2																								
(0.00, 8.00)	1.6886	1.6809	1.6806	1.6439	1.6453	1.6431	1.5827	1.6519	1.6520	1.6531	1.6532	1.6530	1.6532	1.6520	1.6521	1.6532	1.6533	1.6532	1.6533	1.6515	1.6573	1.6574	0.0071	
(0.00, 0.20)	4.4315	4.2497	4.2438	3.4542	3.4732	3.4421	2.3241	4.0054	4.0323	4.0364	4.0711	4.4816	4.6814	4.0078	4.0334	4.0369	4.0724	4.4914	4.6828	3.8731	4.1568	4.1996	0.3267	
(0.75, 8.00)	3.9254	3.9176	3.9173	3.8746	3.8761	3.8767	3.7711	3.8857	3.8858	3.8876	3.8878	3.8907	3.8907	3.8858	3.8860	3.8878	3.8879	3.8907	3.8909	3.8855	3.8920	3.8922	0.0862	
(0.75, 0.20)	10.3822	10.1033	10.0933	8.7123	8.7632	8.7791	6.1037	9.7274	9.7690	9.7836	9.8271	10.3951	10.5815	9.7290	9.7722	9.7861	9.8327	10.4143	10.5924	9.4696	9.9706	10.0230	0.3366	
Case 3																								
(0.00, 8.00)	0.9468	0.9419	0.9417	0.9199	0.9206	0.9196	0.8831	0.9235	0.9235	0.9241	0.9241	0.9241	0.9242	0.9235	0.9236	0.9241	0.9242	0.9242	0.9243	0.9234	0.9269	0.9269	0.0325	
(0.00, 0.20)	3.6449	3.5396	3.5374	3.1108	3.1236	3.1039	2.4700	3.3443	3.3497	3.3640	3.3703	3.4094	3.4301	3.3451	3.3503	3.3658	3.3707	3.4116	3.4325	3.3055	3.4203	3.4290	0.5790	
(0.75, 8.00)	2.3412	2.3368	2.3365	2.3135	2.3140	2.3153	2.2735	2.3177	2.3177	2.3183	2.3183	2.3204	2.3205	2.3177	2.3177	2.3183	2.3183	2.3206	2.3207	2.3175	2.3206	2.3208	0.0756	
(0.75, 0.20)	9.1638	9.0066	9.0002	8.1883	8.2155	8.2605	6.8653	8.6568	8.6717	8.6884	8.7019	8.9460	8.9677	8.6584	8.6731	8.6909	8.7054	8.9505	8.9721	8.5775	8.8133	8.8269	1.1477	

Table 9: MISE: Model 1 (low-dimensional)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058	0.0058
(0.00, 0.20)	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0058	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023	0.0023
(0.75, 8.00)	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030
(0.75, 0.20)	0.0025	0.0026	0.0026	0.0032	0.0032	0.0035	0.0056	0.0028	0.0027	0.0028	0.0027	0.0031	0.0030	0.0028	0.0027	0.0027	0.0027	0.0031	0.0030	0.0029	0.0027	0.0027	0.0111
Case 2																							
(0.00, 8.00)	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140	0.0140
(0.00, 0.20)	0.1301	0.1254	0.1253	0.1033	0.1040	0.1074	0.0666	0.1185	0.1205	0.1193	0.1216	0.1391	0.1402	0.1185	0.1206	0.1193	0.1216	0.1393	0.1402	0.1160	0.1227	0.1247	0.0213
(0.75, 8.00)	0.3377	0.3366	0.3365	0.3300	0.3302	0.3299	0.3163	0.3315	0.3315	0.3317	0.3317	0.3316	0.3316	0.3315	0.3315	0.3317	0.3317	0.3316	0.3316	0.3315	0.3325	0.3325	0.0223
(0.75, 0.20)	0.9133	0.8864	0.8855	0.7625	0.7664	0.7621	0.5452	0.8516	0.8549	0.8574	0.8614	0.9000	0.9139	0.8518	0.8552	0.8580	0.8616	0.9002	0.9142	0.8296	0.8733	0.8771	0.0735
Case 3																							
(0.00, 8.00)	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113	0.0113
(0.00, 0.20)	0.1196	0.1166	0.1164	0.1018	0.1022	0.1037	0.0782	0.1104	0.1108	0.1109	0.1114	0.1169	0.1184	0.1104	0.1109	0.1109	0.1114	0.1169	0.1185	0.1092	0.1129	0.1134	0.0173
(0.75, 8.00)	0.2014	0.2008	0.2008	0.1968	0.1969	0.1974	0.1882	0.1976	0.1976	0.1977	0.1977	0.1984	0.1984	0.1976	0.1976	0.1977	0.1977	0.1984	0.1984	0.1976	0.1983	0.1983	0.0385
(0.75, 0.20)	0.8051	0.7906	0.7900	0.7198	0.7219	0.7253	0.6038	0.7605	0.7617	0.7637	0.7645	0.7831	0.7851	0.7607	0.7620	0.7640	0.7647	0.7834	0.7854	0.7539	0.7721	0.7740	0.1217

Table 10: MISE: Model 2 (high-dimensional, small equal coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ [*]	ESC ₂ [*]	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ [*]	EgMDL ₂ [*]	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01
(0.00, 0.20)	0.06	0.06	0.06	0.09	0.09	0.10	0.15	0.07	0.07	0.07	0.07	0.08	0.08	0.07	0.07	0.07	0.07	0.08	0.08	0.07	0.06	0.06	0.32
(0.75, 8.00)	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.02
(0.75, 0.20)	0.12	0.11	0.11	0.13	0.14	0.17	0.26	0.11	0.11	0.11	0.11	0.14	0.14	0.11	0.11	0.11	0.11	0.14	0.14	0.12	0.11	0.11	0.80
Case 2																							
(0.00, 8.00)	8.06	8.02	8.02	7.84	7.85	7.90	7.47	7.89	7.89	7.90	7.90	7.97	7.97	7.89	7.89	7.90	7.90	7.97	7.97	7.89	7.92	7.92	0.69
(0.00, 0.20)	20.45	19.63	19.59	15.96	16.06	15.96	10.53	18.52	18.68	18.66	18.83	21.21	22.03	18.53	18.68	18.68	18.85	21.31	22.09	17.92	19.22	19.38	1.56
(0.75, 8.00)	27.42	27.35	27.34	26.98	26.99	26.96	26.18	27.06	27.06	27.08	27.08	27.06	27.06	27.06	27.06	27.08	27.08	27.06	27.06	27.06	27.12	27.12	1.53
(0.75, 0.20)	74.17	72.11	72.05	62.18	62.56	62.45	44.15	69.28	69.54	69.80	70.15	73.86	74.90	69.31	69.56	69.82	70.16	74.01	74.93	67.70	71.11	71.43	5.06
Case 3																							
(0.00, 8.00)	5.00	4.98	4.98	4.87	4.87	4.89	4.65	4.89	4.89	4.89	4.89	4.93	4.93	4.89	4.89	4.89	4.89	4.93	4.93	4.89	4.91	4.91	0.72
(0.00, 0.20)	19.32	18.83	18.81	16.72	16.77	16.79	13.49	17.85	17.89	17.95	17.98	18.28	18.36	17.86	17.90	17.96	17.99	18.29	18.38	17.68	18.22	18.26	3.15
(0.75, 8.00)	16.88	16.85	16.85	16.66	16.66	16.68	16.29	16.70	16.70	16.70	16.70	16.73	16.73	16.70	16.70	16.71	16.71	16.73	16.73	16.70	16.72	16.72	1.54
(0.75, 0.20)	67.25	66.10	66.06	60.64	60.81	61.11	51.85	63.80	63.90	64.00	64.11	65.62	65.80	63.84	63.92	64.03	64.12	65.66	65.82	63.32	64.80	64.90	7.78

Table 11: MISE: Model 3 (high-dimensional, decaying coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ [*]	ESC ₂ [*]	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ [*]	EgMDL ₂ [*]	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.66	0.67
(0.00, 0.20)	0.31	0.35	0.35	0.76	0.74	0.81	1.87	0.39	0.37	0.39	0.37	0.29	0.29	0.39	0.37	0.38	0.37	0.29	0.29	0.48	0.34	0.33	1.28
(0.75, 8.00)	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29
(0.75, 0.20)	0.93	0.95	0.95	1.19	1.18	1.35	2.26	1.02	1.00	1.01	1.00	1.15	1.13	1.01	1.00	1.01	1.00	1.14	1.13	1.07	0.98	0.97	4.68
Case 2																							
(0.00, 8.00)	19.76	19.76	19.76	19.75	19.75	19.76	19.71	19.75	19.75	19.75	19.75	19.76	19.76	19.75	19.75	19.75	19.75	19.76	19.76	19.75	19.75	19.75	4.04
(0.00, 0.20)	68.45	65.90	65.78	54.12	54.42	55.41	35.22	62.20	62.96	62.68	63.59	72.29	73.54	62.20	63.00	62.71	63.61	72.37	73.54	60.61	64.34	65.31	8.76
(0.75, 8.00)	146.95	146.46	146.45	143.74	143.81	143.58	138.52	144.31	144.32	144.39	144.40	144.26	144.26	144.32	144.32	144.40	144.41	144.27	144.28	144.30	144.71	144.71	16.65
(0.75, 0.20)	393.43	381.72	381.36	328.81	330.24	329.74	230.64	366.51	368.16	368.76	370.43	391.64	396.30	366.64	368.35	368.96	370.55	392.28	396.51	357.07	376.77	378.52	40.41
Case 3																							
(0.00, 8.00)	13.22	13.21	13.21	13.19	13.19	13.21	13.06	13.19	13.19	13.19	13.19	13.21	13.21	13.19	13.19	13.19	13.19	13.21	13.21	13.19	13.20	13.20	4.38
(0.00, 0.20)	60.78	59.36	59.30	52.64	52.85	53.18	42.02	56.29	56.46	56.52	56.69	58.33	58.66	56.32	56.48	56.57	56.70	58.37	58.68	55.73	57.50	57.66	9.78
(0.75, 8.00)	86.71	86.44	86.43	84.91	84.94	85.08	81.85	85.20	85.20	85.24	85.24	85.44	85.44	85.20	85.20	85.24	85.24	85.44	85.44	85.19	85.43	85.44	15.93
(0.75, 0.20)	360.58	353.80	353.61	322.41	323.57	325.97	270.20	341.00	341.60	342.28	343.03	351.54	352.63	341.07	341.64	342.33	343.11	351.62	352.73	338.38	346.16	346.85	51.89

Table 12: MISE: Model 4 (high-dimensional, slowly decaying coefficients)

1.2 Median number of iterations

The MISE values reported in Tables 1-12 are computed from $N_{TR} = 100$ trials. For each trial, we record the number of iterations that correspond to the models deemed to be optimal according to the stopping rules employed in the experiments. For each quintuple (Case, Model, ω, ζ^2, n) and for each stopping rule, we compute the median of these records. The results are shown in Tables 13-24. In these tables, the convention for each row is to represent in red the smallest value and in bold the second smallest value.

1.2.1 Experiments for $n = 20$

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	420.0	1267.5	381.5	1267.5	1267.5	1267.5	1267.0	803.0	803.0	701.5	741.5	43.5	41.0	762.0	765.5	701.5	701.5	48.5	51.5	765.5	1239.0	1239.0	300.5
(0.00, 0.20)	445.0	1273.5	383.5	1273.5	1273.5	1273.5	1264.5	1089.0	1089.0	1066.0	1066.0	1.0	1.0	1076.0	1076.0	1031.5	1066.0	1.0	1.0	1076.0	1256.0	1256.0	985.0
(0.75, 8.00)	348.0	2105.0	247.0	2105.0	2105.0	2105.0	2024.0	1377.0	1377.0	526.5	526.5	23.0	23.0	1194.5	1377.0	493.0	526.5	24.0	24.0	1377.0	2103.5	2103.5	40.0
(0.75, 0.20)	358.5	2208.0	246.0	2208.0	2208.0	2208.0	2069.5	2063.0	2063.0	2063.0	2063.0	1.0	1.0	2063.0	2063.0	2063.0	2063.0	1.0	1.0	2063.0	2207.5	2207.5	1754.0
Case 2																							
(0.00, 8.00)	834.5	5101.5	489.0	5101.5	5101.5	5101.5	4697.0	2139.5	2139.5	1431.5	1552.0	27.5	28.0	2139.5	2139.5	1552.0	1552.0	28.5	29.5	2139.5	5086.0	5086.0	32.5
(0.00, 0.20)	943.0	5289.0	503.5	5289.0	5289.0	5289.0	5218.5	4653.5	4653.5	4286.5	4286.5	21.0	21.0	4511.0	4511.0	4266.5	4266.5	22.0	22.0	4511.0	5228.0	5228.0	34.5
(0.75, 8.00)	778.0	8134.0	251.0	8134.0	8134.0	8134.0	37.0	128.5	128.5	78.0	80.5	36.0	36.0	139.0	139.0	82.5	86.0	37.0	37.0	139.0	7859.0	7859.0	17.0
(0.75, 0.20)	683.5	8449.0	253.5	8449.0	8449.0	8449.0	21.0	6386.5	6386.5	4435.0	4435.0	21.0	21.0	6224.0	6224.0	4400.5	4400.5	22.0	22.0	6224.0	8426.0	8426.0	24.0
Case 3																							
(0.00, 8.00)	583.0	2600.0	383.0	2600.0	2600.0	2600.0	2513.0	1110.5	1249.0	866.0	866.0	36.5	37.0	1021.5	1021.5	866.0	866.0	38.5	38.5	1055.0	2588.0	2588.0	75.5
(0.00, 0.20)	543.0	2768.0	329.5	2768.0	2768.0	2768.0	2615.5	2352.5	2352.5	2058.5	2058.5	18.0	17.5	2174.0	2237.5	2031.5	2031.5	21.0	21.0	2237.5	2755.5	2755.5	130.5
(0.75, 8.00)	464.0	4366.5	174.0	4366.5	4366.5	4366.0	1256.0	474.5	474.5	82.5	93.0	29.5	29.5	474.5	474.5	93.0	110.0	30.0	30.0	474.5	4360.0	4360.0	31.0
(0.75, 0.20)	458.5	4442.0	221.0	4442.0	4442.0	4442.0	26.5	3752.5	3793.5	3576.0	3592.5	12.0	11.5	3665.5	3740.5	3576.0	3576.0	14.5	13.0	3740.5	4442.0	4442.0	69.0

Table 13: Median number of iterations: Model 1 (low-dimensional)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	413.0	1244.5	340.0	1244.5	1244.5	1244.5	1215.0	876.0	882.0	831.0	831.0	1.0	1.0	835.5	835.5	831.0	831.0	1.0	5.0	835.5	1238.5	1238.5	952.0
(0.00, 0.20)	405.0	1285.5	345.5	1285.5	1285.5	1285.5	1279.0	1100.0	1100.0	1082.5	1090.0	1.0	1.0	1090.0	1090.0	1082.5	1082.5	1.0	1.0	1090.0	1283.5	1283.5	1040.0
(0.75, 8.00)	294.0	1891.5	222.0	1891.5	1891.5	1891.5	1622.0	782.0	807.0	591.5	613.0	26.5	23.5	665.0	665.0	581.5	609.5	34.5	39.5	665.0	1890.5	1890.5	1250.5
(0.75, 0.20)	329.0	2210.0	215.5	2210.0	2210.0	2210.0	1718.0	2039.0	2039.0	2026.5	2026.5	1.0	1.0	2029.0	2029.0	2006.5	2006.5	1.0	1.0	2029.0	2185.5	2185.5	1611.5
Case 2																							
(0.00, 8.00)	810.0	4984.0	297.5	4984.0	4984.0	4984.0	4817.0	4284.5	4284.5	3400.5	3400.5	25.0	25.0	4014.0	4014.0	2798.0	2798.0	26.0	26.0	4014.0	4983.0	4983.0	60.0
(0.00, 0.20)	808.0	5159.5	346.0	5159.5	5159.5	5159.5	4809.5	4523.5	4532.0	4014.0	4062.0	17.5	19.5	4488.0	4488.0	3918.0	3951.0	20.0	20.0	4488.0	5159.5	5159.5	63.0
(0.75, 8.00)	844.0	8288.0	288.0	8288.0	8288.0	8287.5	30.0	1092.5	1093.0	315.0	373.5	29.0	29.0	1092.5	1093.0	265.5	337.5	30.0	30.0	1093.0	8053.5	8053.5	37.0
(0.75, 0.20)	677.5	8413.0	185.5	8413.0	8413.0	8413.0	19.0	5607.5	5742.0	3219.0	4405.0	20.0	20.0	5607.5	5607.5	3491.5	4405.0	21.0	21.5	5607.5	8390.5	8390.5	35.5
Case 3																							
(0.00, 8.00)	549.5	2538.0	302.0	2538.0	2538.0	2538.0	2447.0	2062.5	2062.5	1754.0	1813.5	18.0	18.0	2002.5	2002.5	1725.5	1742.5	21.0	21.0	2002.5	2531.0	2531.0	1001.0
(0.00, 0.20)	481.0	2606.0	340.5	2606.0	2606.0	2606.0	2532.0	2277.5	2277.5	2154.5	2154.5	17.0	16.5	2277.5	2277.5	2116.0	2116.0	17.0	17.0	2277.5	2593.5	2593.5	160.5
(0.75, 8.00)	505.5	3966.5	300.5	3966.5	3966.5	3966.5	3402.5	1671.5	1762.0	890.0	959.5	28.5	28.5	1671.5	1671.5	857.5	907.5	31.0	31.0	1671.5	3882.0	3882.0	259.5
(0.75, 0.20)	459.5	4175.0	223.0	4175.0	4175.0	4175.0	6.5	3247.5	3269.0	3115.0	3115.0	12.0	10.5	3123.5	3159.0	2932.0	3032.5	13.5	13.0	3159.0	4162.0	4162.0	122.5

Table 14: Median number of iterations: Model 2 (high-dimensional, small equal coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	401.5	1199.0	347.5	1199.0	1199.0	1199.0	1192.5	812.5	812.5	696.5	769.0	20.5	16.0	807.0	807.0	687.5	696.5	24.5	22.5	807.0	1190.0	1190.0	65.5
(0.00, 0.20)	466.0	1318.5	358.5	1318.5	1318.5	1318.5	1318.5	1199.0	1199.0	1068.5	1068.5	1.0	1.0	1073.5	1073.5	1022.5	1022.5	1.0	3.0	1073.5	1309.0	1309.0	1080.5
(0.75, 8.00)	331.0	1971.0	238.5	1971.0	1971.0	1971.0	1859.5	837.0	845.5	523.0	548.0	32.0	32.0	778.0	813.5	491.5	535.0	33.0	33.0	813.5	1964.0	1964.0	196.0
(0.75, 0.20)	366.5	2110.0	241.5	2110.0	2110.0	2110.0	1719.0	2028.5	2028.5	1999.0	1999.0	1.0	1.0	2006.5	2028.5	1960.5	1974.5	1.0	1.0	2028.5	2099.0	2099.0	1321.5
Case 2																							
(0.00, 8.00)	737.5	5169.5	412.5	5169.5	5169.5	5169.5	4595.0	1536.0	1536.0	793.5	853.0	31.0	31.0	1329.0	1326.0	767.0	793.5	32.0	32.0	1483.0	5076.0	5076.0	16.0
(0.00, 0.20)	802.0	5306.0	407.0	5306.0	5306.0	5306.0	4680.0	4321.5	4321.5	3218.0	3335.0	23.0	23.0	4022.5	4022.5	3065.0	3087.5	23.0	23.5	4022.5	5274.0	5274.0	28.5
(0.75, 8.00)	697.5	8447.0	359.5	8447.0	8447.0	8447.0	6265.0	308.5	337.5	97.5	97.5	33.5	33.5	337.5	356.5	99.5	100.5	34.0	34.0	356.5	8355.0	8355.0	16.0
(0.75, 0.20)	652.0	9106.5	217.0	9106.5	9106.5	9106.5	20.5	7408.0	7577.5	4000.0	4000.0	21.0	21.0	6641.5	6950.0	4000.0	4000.0	22.5	22.0	7146.0	9106.0	9106.0	84.5
Case 3																							
(0.00, 8.00)	536.0	2645.0	347.5	2645.0	2645.0	2645.0	2609.0	1819.5	1819.5	1309.5	1384.5	26.0	24.0	1694.0	1735.0	1309.5	1348.5	29.0	29.0	1735.0	2623.0	2623.0	97.5
(0.00, 0.20)	442.0	2680.0	281.5	2680.0	2680.0	2680.0	2424.5	2200.0	2200.0	1958.5	2060.0	14.0	14.0	2105.5	2128.5	1884.5	1884.5	15.5	15.5	2168.5	2669.5	2669.5	176.0
(0.75, 8.00)	493.5	4182.5	233.5	4182.5	4182.5	4182.5	3297.0	446.0	488.0	175.5	195.5	37.0	37.0	446.0	455.5	193.0	195.5	38.0	38.0	455.5	4165.5	4165.5	66.5
(0.75, 0.20)	426.5	4423.5	178.0	4423.5	4423.5	4423.5	12.5	3877.5	3939.0	3603.5	3655.5	15.0	14.0	3771.5	3771.5	2885.0	3225.0	17.0	17.0	3851.5	4405.5	4405.5	64.0

Table 15: Median number of iterations: Model 3 (high-dimensional, decaying coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	397.5	1203.0	333.0	1203.0	1203.0	1203.0	1199.5	1110.5	1110.5	990.0	1055.0	1.0	1.0	1104.5	1110.5	932.0	943.5	1.0	6.0	1110.5	1194.5	1194.5	887.5
(0.00, 0.20)	428.5	1292.5	361.0	1292.5	1292.5	1292.5	1259.5	1088.5	1088.5	1072.5	1072.5	1.0	1.0	1074.5	1074.5	995.5	1045.0	1.0	1.5	1074.5	1277.5	1277.5	907.5
(0.75, 8.00)	324.0	1905.5	234.5	1905.5	1905.5	1905.5	1787.0	876.0	966.0	512.0	633.5	31.0	31.5	819.0	876.0	428.5	570.0	36.5	39.0	876.0	1881.5	1881.5	1141.5
(0.75, 0.20)	321.0	2145.5	183.5	2145.5	2145.5	2145.5	1.5	2017.5	2056.0	2000.0	2000.0	1.0	1.0	2017.0	2017.0	1999.5	1999.5	1.0	2.0	2017.0	2132.0	2132.0	135.0
Case 2																							
(0.00, 8.00)	771.5	4994.0	443.0	4994.0	4994.0	4994.0	4539.5	3406.5	3406.5	3092.5	3183.5	22.5	23.0	3406.5	3406.5	2993.5	3092.5	23.0	23.5	3406.5	4893.0	4893.0	106.0
(0.00, 0.20)	803.0	4986.0	420.0	4986.0	4986.0	4986.0	4656.0	3428.0	3556.5	2730.0	2730.0	25.0	25.0	3428.0	3525.5	2537.5	2730.0	25.5	26.0	3525.5	4969.5	4969.5	23.5
(0.75, 8.00)	773.5	8495.0	397.0	8495.0	8495.0	8495.0	5996.0	1343.5	1480.0	236.5	295.0	30.0	30.0	1044.5	1051.0	295.0	310.5	31.5	32.0	1053.0	8391.5	8391.5	24.5
(0.75, 0.20)	763.0	8618.5	201.0	8618.5	8618.5	8618.5	26.5	4861.5	4861.5	4659.0	4659.0	22.0	22.0	4861.5	4861.5	4659.0	4659.0	22.5	23.0	4861.5	8485.0	8485.0	40.5
Case 3																							
(0.00, 8.00)	490.5	2566.0	378.0	2566.0	2566.0	2566.0	2357.0	1768.0	1800.5	1635.0	1635.0	18.0	15.5	1746.0	1746.0	1528.0	1528.0	19.0	18.5	1746.0	2564.5	2564.5	200.5
(0.00, 0.20)	520.0	2611.0	287.5	2611.0	2611.0	2611.0	2456.0	2282.5	2282.5	2174.0	2174.0	16.5	16.0	2228.5	2228.5	2120.0	2120.0	19.5	20.0	2228.5	2609.5	2609.5	438.5
(0.75, 8.00)	588.0	4424.0	280.0	4424.0	4424.0	4424.0	3463.0	1643.0	1643.0	679.0	783.0	27.5	26.5	1287.0	1585.5	659.5	753.5	29.5	29.5	1585.5	4424.0	4424.0	227.0
(0.75, 0.20)	416.0	4612.0	156.0	4612.0	4612.0	4612.0	12.0	4168.0	4168.0	4103.5	4103.5	11.0	8.0	4168.0	4168.0	4103.5	4103.5	14.0	13.5	4168.0	4593.5	4593.5	101.5

Table 16: Median number of iterations: Model 4 (high-dimensional, slowly decaying coefficients)

1.2.2 Experiments for $n = 100$

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	423.5	19985.5	248.0	188.0	164.5	75.0	59.0	144.5	136.5	109.0	106.0	71.0	71.0	151.0	146.5	114.5	112.0	72.0	72.0	145.0	202.5	191.5	96.0
(0.00, 0.20)	352.5	19982.5	177.5	115.5	70.0	5.0	1.0	564.5	538.0	617.0	568.5	5.0	1.0	573.5	577.5	649.5	605.5	9.0	7.0	594.0	17368.0	19979.0	26.5
(0.75, 8.00)	478.0	19990.5	136.0	85.0	62.0	35.0	29.0	58.5	57.5	47.0	46.5	33.0	33.0	62.0	62.0	49.0	49.0	33.0	33.0	62.0	57.5	56.0	46.0
(0.75, 0.20)	428.0	19970.5	105.5	54.5	35.0	8.5	3.0	1954.5	1368.5	3468.5	3373.0	8.5	8.0	1989.5	1401.0	3797.0	4082.0	10.0	9.5	2030.5	19970.5	19999.0	18.5
Case 2																							
(0.00, 8.00)	5042.5	19997.0	3296.0	4422.0	5848.0	861.0	124.0	1043.0	1043.0	963.5	987.5	367.5	367.5	1043.0	1043.0	988.5	988.5	378.5	378.5	1043.0	1386.5	1386.5	7.0
(0.00, 0.20)	4002.5	19984.0	1999.0	2304.0	2893.5	140.0	39.0	1156.0	1151.5	1013.5	1028.5	98.5	106.0	1244.0	1262.5	1185.0	1185.5	106.0	108.5	1262.5	2725.0	2725.0	5.0
(0.75, 8.00)	9363.0	19999.0	3915.5	5152.0	6737.5	265.0	61.5	631.0	631.0	491.5	507.5	116.5	116.5	648.0	648.0	534.5	534.5	117.5	119.0	648.0	919.0	950.5	6.0
(0.75, 0.20)	8635.5	19999.0	3276.5	4453.5	5317.5	132.0	37.0	1444.0	1444.0	1381.0	1346.5	111.0	112.0	1500.5	1549.0	1432.5	1452.5	118.0	118.0	1500.5	5401.5	5403.5	5.0
Case 3																							
(0.00, 8.00)	1299.0	19971.5	753.5	752.5	761.5	191.0	94.5	324.0	326.5	276.0	279.0	152.0	155.0	331.0	329.5	284.0	284.0	155.5	155.5	329.5	413.0	412.5	41.5
(0.00, 0.20)	1103.5	19699.5	507.0	421.5	400.0	73.0	22.0	576.0	574.5	532.5	534.5	71.0	70.0	599.0	589.5	546.0	574.5	74.5	73.5	581.0	1944.0	1915.5	35.5
(0.75, 8.00)	2626.0	19988.0	881.5	833.0	694.5	91.0	45.0	172.0	172.0	140.5	140.5	71.5	71.0	178.0	178.0	151.5	151.5	72.5	73.0	178.0	279.5	279.5	14.5
(0.75, 0.20)	2648.5	19986.5	833.0	547.5	365.5	54.0	20.5	1221.0	1186.5	1146.5	929.5	50.5	51.0	1305.0	1301.5	1214.5	1186.5	56.0	52.5	1214.0	8691.5	8968.5	32.0

Table 17: Median number of iterations: Model 1 (low-dimensional)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	903.0	19994.0	541.5	14873.5	16885.0	17388.0	1.0	844.5	992.5	974.0	1229.5	64.0	27.5	848.5	992.5	974.0	1229.5	84.0	37.0	1046.0	2015.0	2697.5	3295.0
(0.00, 0.20)	346.5	19995.0	197.0	121.5	83.5	2.0	1.0	536.0	531.5	571.5	596.5	1.0	1.0	550.5	532.5	592.0	603.5	1.0	3.0	570.5	17263.5	19983.5	37.0
(0.75, 8.00)	636.0	19975.5	389.0	383.5	391.0	252.0	32.0	307.5	318.0	309.5	324.5	212.0	221.0	308.5	322.5	309.5	329.5	212.0	224.0	322.5	394.5	420.0	1016.0
(0.75, 0.20)	355.0	19978.5	120.5	79.0	58.5	10.0	1.0	1725.5	939.5	3136.0	2782.0	12.0	5.0	1759.5	1213.0	3394.5	3019.0	17.0	11.5	746.0	19971.5	20000.0	48.5
Case 2																							
(0.00, 8.00)	12068.0	20000.0	6234.0	19999.0	19999.0	19997.5	44.0	2935.5	2935.5	3545.5	3544.5	204.0	208.0	2954.5	2954.5	3707.0	3706.0	222.5	227.0	2954.5	6473.0	6475.5	5.0
(0.00, 0.20)	4384.5	19997.0	1922.0	3007.0	4237.5	137.0	37.0	1282.5	1302.0	1282.5	1302.0	106.5	106.5	1303.5	1325.0	1329.0	1336.0	108.0	109.5	1326.5	3204.5	3001.0	5.0
(0.75, 8.00)	14151.0	20000.0	8195.5	16812.0	18733.0	4335.0	59.0	3034.0	2995.0	3096.5	3038.5	829.5	829.5	3034.0	3035.0	3096.5	3096.5	829.5	829.5	3035.0	4246.0	4258.0	6.0
(0.75, 0.20)	8801.5	19998.0	3594.0	5339.5	6600.5	120.5	37.0	1881.5	1897.0	1573.5	1494.5	95.0	95.5	2096.5	2103.0	1832.5	1728.5	97.0	99.0	2103.0	5793.0	5797.0	6.0
Case 3																							
(0.00, 8.00)	3527.5	20000.0	2056.0	19980.5	19990.5	19997.0	26.0	1743.5	1804.5	2162.5	2218.0	147.5	149.5	1764.0	1820.0	2162.5	2218.0	151.0	150.5	1905.5	3500.0	3713.5	199.0
(0.00, 0.20)	1199.0	19437.0	606.5	585.5	576.0	62.5	24.0	635.5	631.5	599.0	615.5	63.0	60.5	661.0	665.5	670.0	644.5	66.5	63.5	651.5	1854.5	1941.0	39.0
(0.75, 8.00)	3820.5	19998.0	1842.5	3181.0	5136.5	708.5	51.5	921.0	949.0	941.5	980.0	373.0	373.0	949.0	949.5	979.5	980.0	373.0	377.5	949.5	1329.5	1349.0	148.0
(0.75, 0.20)	2495.0	19993.5	820.5	501.0	376.0	54.0	23.0	947.5	879.5	989.0	831.0	54.5	54.0	1050.0	962.5	989.0	941.0	56.5	56.0	962.5	8947.5	9076.0	53.5

Table 18: Median number of iterations: Model 2 (high-dimensional, small equal coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ [*]	ESC ₂ [*]	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ [*]	EgMDL ₂ [*]	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	560.0	19985.0	372.5	389.0	383.5	105.5	37.0	303.0	301.5	282.5	281.0	95.0	92.5	309.5	306.0	284.0	284.0	99.5	94.5	306.0	466.5	463.0	279.5
(0.00, 0.20)	329.0	17492.5	172.0	98.5	59.5	5.5	1.0	589.0	502.0	641.5	567.5	6.0	1.0	593.0	502.0	642.0	594.0	10.0	7.5	510.5	14494.5	19623.0	24.0
(0.75, 8.00)	490.0	19967.0	231.0	180.0	169.0	91.5	44.0	148.5	149.0	142.0	142.0	82.0	82.0	150.5	151.5	147.0	147.0	83.0	83.0	151.5	182.5	184.0	198.5
(0.75, 0.20)	477.5	19987.5	134.5	60.0	38.5	9.0	2.0	1662.0	1372.5	3482.5	3540.5	10.0	8.0	1680.5	1533.0	3570.5	4145.5	12.0	10.5	1552.0	19979.0	20000.0	25.0
Case 2																							
(0.00, 8.00)	6887.5	19997.5	4148.5	10257.0	14111.0	2045.5	79.5	1673.5	1674.0	1649.0	1649.5	389.0	389.0	1702.5	1703.0	1783.0	1783.0	393.5	404.0	1716.0	2695.0	2712.0	6.0
(0.00, 0.20)	3650.0	19993.5	1930.5	2601.0	2797.5	114.5	39.0	1054.5	1058.5	1054.5	1045.0	96.0	96.0	1118.5	1096.0	1115.0	1058.5	100.5	104.0	1119.5	2362.0	2350.5	5.0
(0.75, 8.00)	12158.5	20000.0	4980.0	11507.0	13987.5	813.0	82.0	1186.5	1180.0	1055.0	1055.0	263.0	263.0	1227.5	1227.5	1055.0	1115.5	263.0	263.0	1227.5	1988.0	1988.0	5.0
(0.75, 0.20)	8727.5	19999.0	3797.0	5763.5	7207.5	149.5	35.0	2296.0	2321.0	2348.5	2349.0	113.5	113.5	2385.0	2321.0	2348.5	2384.5	118.5	121.0	2349.0	7518.0	7031.0	5.0
Case 3																							
(0.00, 8.00)	1975.5	19997.0	1252.0	2320.5	3283.0	318.0	61.0	593.5	599.5	589.5	589.5	175.0	172.5	607.0	630.5	599.5	615.5	181.0	185.0	630.5	1052.0	1051.0	51.0
(0.00, 0.20)	1054.5	19955.5	552.5	463.5	436.5	57.0	21.0	630.0	630.0	577.0	561.0	58.0	56.5	642.5	650.0	589.5	582.0	63.0	57.5	630.0	1901.5	1850.0	40.0
(0.75, 8.00)	2606.5	19993.5	1194.0	1310.5	1310.5	152.0	62.0	311.5	309.5	261.5	264.0	111.0	112.0	323.0	323.0	271.0	279.5	117.0	117.5	323.0	553.0	560.5	21.5
(0.75, 0.20)	2175.0	19995.5	875.5	665.5	542.5	55.5	22.0	874.0	887.0	788.5	800.0	54.0	52.0	951.0	951.0	803.0	814.0	56.5	55.5	885.0	6407.0	5900.5	30.5

Table 19: Median number of iterations: Model 3 (high-dimensional, decaying coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ [*]	ESC ₂ [*]	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ [*]	EgMDL ₂ [*]	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	900.5	19987.0	562.0	6555.5	13479.5	16160.0	8.0	733.0	800.0	804.5	890.0	112.5	84.5	735.0	804.0	804.5	899.5	115.5	96.0	840.0	1347.5	1527.0	2226.5
(0.00, 0.20)	349.0	19962.0	181.5	108.0	76.5	2.0	1.0	561.5	541.0	657.0	657.5	1.0	1.0	572.0	569.0	686.5	699.0	1.0	3.5	604.5	16456.5	19939.5	17.5
(0.75, 8.00)	560.5	19973.5	358.5	335.0	335.5	217.0	43.0	288.0	293.5	288.0	293.5	183.0	188.5	289.5	294.0	289.5	293.5	187.5	190.0	294.0	342.0	355.5	559.5
(0.75, 0.20)	449.5	19892.5	126.0	75.0	54.0	8.0	1.0	1355.5	978.0	2571.5	1485.5	10.0	5.0	1355.5	1024.0	2715.0	2695.0	13.0	9.5	1019.0	19979.5	20000.0	42.5
Case 2																							
(0.00, 8.00)	10147.0	20000.0	6187.5	19995.0	19995.5	19985.5	48.0	2943.5	2943.5	3561.0	3606.5	299.5	299.5	2943.5	2993.0	3634.5	3682.5	311.0	311.0	3002.0	5305.0	5305.0	6.0
(0.00, 0.20)	3951.0	19991.5	2213.5	3251.0	3522.5	123.0	38.5	1236.0	1193.5	1213.5	1173.0	99.0	101.5	1359.0	1359.5	1310.5	1269.5	102.5	102.5	1396.5	2931.5	2931.5	5.0
(0.75, 8.00)	13292.0	20000.0	6457.0	13718.0	18474.5	2609.0	72.5	2335.5	2335.5	2482.5	2483.0	598.5	598.0	2335.5	2335.5	2482.5	2483.0	598.5	604.0	2335.5	3240.0	3240.0	5.0
(0.75, 0.20)	8396.5	19999.0	4096.0	5241.5	7862.0	97.0	33.5	1920.0	1921.0	1832.0	1790.0	73.5	76.0	2158.5	2055.5	2026.5	1904.0	81.0	79.0	2026.5	7099.0	7000.5	5.0
Case 3																							
(0.00, 8.00)	3013.0	20000.0	1737.0	19344.5	19983.0	19994.5	33.5	1350.0	1358.5	1462.0	1502.0	192.5	199.0	1356.0	1393.5	1493.0	1502.0	217.0	204.0	1393.5	2708.5	2905.5	140.5
(0.00, 0.20)	1243.0	19994.5	551.0	477.0	497.5	59.0	21.5	602.0	616.5	575.0	583.5	59.0	58.5	661.0	636.5	649.5	598.0	62.0	61.5	649.0	2730.5	2683.5	36.5
(0.75, 8.00)	3528.5	19996.0	1690.5	2202.0	3137.5	422.0	61.5	617.5	642.0	640.5	642.0	274.5	274.5	647.0	662.5	647.0	654.0	274.5	274.5	662.5	1079.0	1079.5	58.0
(0.75, 0.20)	1944.5	19995.5	668.5	467.0	315.0	62.5	21.0	796.0	757.0	623.5	592.5	62.0	62.0	909.0	777.5	638.5	640.0	62.5	62.5	757.0	8460.5	8417.5	50.0

Table 20: Median number of iterations: Model 4 (high-dimensional, slowly decaying coefficients)

1.2.3 Experiments for $n = 10000$

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	348.0	258.0	258.0	148.0	139.0	132.0	124.5	143.0	143.0	136.0	136.0	131.0	131.0	144.0	144.0	136.0	136.0	131.0	131.0	143.5	144.0	143.5	401.0
(0.00, 0.20)	295.5	205.0	204.5	93.0	86.0	77.5	71.0	119.5	117.0	102.0	100.0	82.0	82.0	122.5	122.0	105.0	103.5	82.0	82.0	101.0	127.5	126.5	3248.0
(0.75, 8.00)	345.5	140.5	140.5	67.5	64.0	60.0	58.0	65.0	65.0	62.0	62.0	60.0	60.0	65.5	65.5	63.0	63.0	60.0	60.0	65.5	64.0	64.0	228.5
(0.75, 0.20)	326.5	119.5	119.5	46.0	41.0	37.0	35.0	52.0	52.5	45.5	45.0	38.0	38.0	55.0	55.0	47.0	47.0	39.0	39.0	47.0	47.0	47.0	9692.5
Case 2																							
(0.00, 8.00)	2131.5	1829.0	1825.0	1001.0	1015.5	969.0	444.5	1118.0	1118.0	1128.0	1128.0	1128.0	1128.0	1118.0	1118.0	1128.0	1128.0	1128.0	1128.0	1106.0	1205.0	1209.5	11.0
(0.00, 0.20)	2078.5	1766.0	1760.5	928.0	946.0	926.5	385.0	1444.0	1461.0	1462.5	1492.0	1923.5	4579.0	1122.5	1461.0	1462.5	1493.5	1944.5	4579.0	1278.0	1631.0	1691.5	16.0
(0.75, 8.00)	8889.0	7392.5	7347.0	3526.0	3625.5	3672.0	779.5	4146.5	4146.5	4323.5	4323.5	4526.5	4526.5	4146.5	4151.5	4363.5	4363.5	4526.5	4526.5	4146.5	4557.5	4579.0	5.0
(0.75, 0.20)	8897.5	7317.5	7279.0	3505.5	3593.5	3581.0	728.0	5964.5	6123.0	6221.5	6369.5	8555.0	10324.0	5974.0	6132.5	6221.5	6369.5	8628.5	10355.5	5251.5	6804.5	7030.0	4.0
Case 3																							
(0.00, 8.00)	3020.0	2506.5	2503.0	1423.0	1442.0	1417.5	724.0	1555.0	1555.0	1565.0	1565.0	1560.0	1561.5	1555.0	1555.0	1565.0	1565.0	1561.5	1565.0	1555.0	1685.5	1685.5	14.0
(0.00, 0.20)	2953.0	2450.5	2445.0	1356.5	1369.5	1348.5	663.0	1824.0	1833.0	1908.5	1929.0	1987.5	1992.0	1824.0	1833.0	1916.5	1929.0	1987.5	1992.0	1740.0	2006.5	2052.5	25.5
(0.75, 8.00)	14253.5	12142.5	12073.0	6306.5	6418.0	6731.5	2213.5	7031.5	7031.5	7138.0	7138.0	7565.5	7565.5	7031.5	7031.5	7138.0	7138.0	7568.0	7568.0	6950.5	7568.0	7583.5	5.0
(0.75, 0.20)	13900.0	11883.5	11748.5	5998.5	6087.0	6583.5	2057.5	8666.0	8775.0	8940.5	8946.5	11137.0	11340.0	8744.5	8775.0	8946.5	8946.5	11137.0	11446.0	8219.0	9861.0	9863.0	7.5

Table 21: Median number of iterations: Model 1 (low-dimensional)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV		
Case 1																									
(0.00, 8.00)	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	266.0	260.0	
(0.00, 0.20)	765.0	765.0	765.0	764.0	765.0	764.0	274.5	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	765.0	739.5	
(0.75, 8.00)	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	223.0	219.5
(0.75, 0.20)	427.0	391.5	391.5	296.5	299.0	274.0	201.5	347.0	354.5	349.0	358.5	310.0	315.5	347.0	355.5	349.0	359.0	310.0	316.0	327.0	364.0	373.0	4698.0		
Case 2																									
(0.00, 8.00)	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	917.5	148.0	
(0.00, 0.20)	2645.5	2318.0	2306.0	1398.0	1414.5	1520.0	651.5	1950.0	2063.5	1984.5	2119.0	4550.0	4619.0	1950.0	2068.0	1984.5	2119.0	4582.0	4619.0	1841.0	2148.5	2299.5	173.5		
(0.75, 8.00)	6371.5	5354.5	5328.5	2501.0	2587.0	2477.0	889.5	2933.0	2933.0	2942.5	2942.5	2941.0	2941.0	2933.0	2933.0	2943.0	2943.0	2941.0	2941.0	2933.0	3302.5	3303.0	37.0		
(0.75, 0.20)	8344.0	6600.5	6594.5	2952.5	3039.0	2935.5	617.5	5255.0	5454.5	5568.0	5735.5	6714.5	7636.0	5279.0	5454.5	5607.5	5740.0	6714.5	7636.0	4494.0	6023.5	6238.5	32.0		
Case 3																									
(0.00, 8.00)	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	1454.0	26.5	
(0.00, 0.20)	3878.0	3420.5	3381.5	1949.0	2007.0	2082.5	915.5	2683.5	2735.5	2728.0	2772.5	3453.5	3687.5	2683.5	2755.5	2728.0	2772.5	3453.5	3687.5	2556.5	2923.0	3028.0	20.5		
(0.75, 8.00)	14286.0	12112.0	12112.0	6724.5	6865.5	7260.5	2954.5	7489.5	7539.0	7539.0	7539.0	8396.0	8396.0	7539.0	7539.0	7539.0	7539.0	8449.0	8449.0	7361.0	8084.0	8084.0	28.0		
(0.75, 0.20)	13761.5	11382.0	11327.0	5778.5	5944.0	6099.0	1958.0	8261.5	8474.0	8766.5	8795.0	10400.5	10623.5	8308.0	8489.5	8766.5	8795.0	10430.5	10623.5	7861.0	9320.0	9471.5	12.0		

Table 22: Median number of iterations: Model 2 (high-dimensional, small equal coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ ⁺	ESC ₂ ⁺	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ ⁺	EgMDL ₂ ⁺	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	2308.0	2307.0	2305.0	1559.0	1571.5	1793.5	795.5	1771.5	1780.0	1790.0	1790.5	2241.5	2295.5	1776.5	1780.0	1790.0	1795.0	2266.0	2295.5	1764.5	1894.5	1899.0	1877.0
(0.00, 0.20)	665.0	517.0	509.0	224.0	221.5	169.0	105.0	389.0	378.0	386.0	376.0	241.0	241.0	396.0	386.5	388.0	378.0	241.5	241.0	295.5	467.5	465.0	3038.5
(0.75, 8.00)	536.5	504.0	504.0	411.0	413.5	391.5	317.5	420.0	421.0	423.0	423.0	398.0	398.0	421.0	422.0	423.0	424.0	398.0	399.0	420.0	430.0	430.5	1012.5
(0.75, 0.20)	338.0	235.0	235.0	137.0	135.0	110.5	73.0	176.5	177.0	174.0	174.0	130.5	130.0	177.5	178.0	175.0	175.5	131.0	131.0	156.5	192.0	192.0	1831.0
Case 2																							
(0.00, 8.00)	3202.0	2831.0	2798.0	1964.5	1984.0	2157.5	1138.0	2147.5	2156.0	2157.5	2166.5	2493.5	2496.5	2153.0	2156.0	2162.0	2166.5	2493.5	2499.5	2153.0	2254.0	2264.0	4.0
(0.00, 0.20)	2146.0	1821.5	1806.0	994.5	1007.0	993.5	398.5	1493.0	1525.0	1543.0	1577.5	3451.0	4475.5	1502.5	1525.0	1550.0	1577.5	3724.0	4501.5	1355.0	1690.5	1725.0	26.5
(0.75, 8.00)	7473.5	5908.5	5790.0	2439.0	2523.5	2320.0	616.5	2863.0	2863.0	3008.0	3008.0	2818.5	2842.0	2878.0	2878.0	3008.0	3008.0	2842.0	2842.0	2850.0	3242.0	3242.0	5.0
(0.75, 0.20)	8749.0	7214.5	7146.0	3261.5	3362.5	3370.5	699.5	5685.5	5731.5	5790.5	5962.5	8100.5	9118.5	5685.5	5790.0	5831.5	5962.5	8288.5	9118.5	4856.5	6571.5	6798.0	8.0
Case 3																							
(0.00, 8.00)	4828.0	4267.5	4229.5	2883.5	2907.5	3161.5	1718.5	3039.5	3056.0	3093.0	3121.0	3459.0	3459.0	3056.0	3056.0	3118.0	3121.0	3459.0	3459.0	3039.5	3280.0	3288.5	6.0
(0.00, 0.20)	3100.5	2619.0	2619.0	1487.5	1526.0	1526.0	723.5	1970.0	2003.0	2012.5	2017.0	2160.5	2187.0	1984.5	2003.0	2012.5	2018.5	2160.5	2226.5	1897.5	2167.5	2185.0	17.5
(0.75, 8.00)	14484.5	12233.5	12166.0	6406.0	6472.5	6799.5	2346.5	7098.0	7098.0	7246.5	7246.5	7892.5	7954.5	7115.0	7122.0	7295.5	7295.5	7954.5	7954.5	7098.0	7843.0	7844.5	6.0
(0.75, 0.20)	13840.5	11637.0	11630.5	6216.0	6360.5	6585.0	2006.5	8688.0	8748.0	8866.5	8994.0	10359.0	10959.0	8688.0	8748.0	8866.5	8994.0	10450.5	10981.0	8162.5	9633.5	9762.5	8.5

Table 23: Median number of iterations: Model 3 (high-dimensional, decaying coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ ⁺	ESC ₂ ⁺	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ ⁺	EgMDL ₂ ⁺	MMLU	MMLG ₁	MMLG ₂	CV	
Case 1																								
(0.00, 8.00)	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	669.5	661.5
(0.00, 0.20)	1877.0	1567.5	1557.0	699.5	717.5	659.0	199.5	1346.5	1427.5	1367.5	1459.5	2675.5	2675.5	1348.5	1432.0	1373.0	1460.5	2675.5	2675.5	1100.0	1589.0	1736.0	2156.5	
(0.75, 8.00)	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	312.0	294.5
(0.75, 0.20)	396.0	359.0	359.0	264.0	265.0	237.0	169.0	312.5	318.5	315.0	321.0	274.0	275.5	314.0	320.0	316.0	322.0	274.0	276.5	291.5	330.0	336.5	7314.5	
Case 2																								
(0.00, 8.00)	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.0	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	1721.5	20.5
(0.00, 0.20)	2446.0	2119.0	2112.5	1237.5	1259.5	1301.5	553.5	1771.5	1850.5	1807.5	1897.0	3928.5	4234.0	1771.5	1853.0	1807.5	1900.0	3997.5	4234.0	1636.5	1980.0	2078.0	98.5	
(0.75, 8.00)	6685.0	5298.0	5261.0	2279.5	2317.0	2215.0	807.0	2672.0	2672.0	2684.5	2684.5	2638.5	2647.0	2672.0	2672.0	2684.5	2719.0	2647.0	2667.5	2936.0	2936.0	2936.5	11.0	
(0.75, 0.20)	8420.5	6939.5	6871.0	3027.0	3105.0	3070.5	676.0	5358.5	5444.0	5480.5	5689.5	7551.5	8290.0	5358.5	5482.5	5532.0	5696.5	7623.5	8290.0	4596.0	6355.0	6483.0	24.0	
Case 3																								
(0.00, 8.00)	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3362.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	3363.5	22.0
(0.00, 0.20)	3539.0	3058.0	3046.5	1755.0	1796.5	1839.5	869.5	2337.5	2354.5	2363.5	2401.0	2782.0	2819.0	2346.0	2360.5	2375.0	2401.0	2788.0	2845.5	2245.0	2593.5	2616.0	18.5	
(0.75, 8.00)	14947.5	12628.5	12536.5	6861.5	6960.0	7209.0	2656.5	7641.0	7641.0	7666.5	7666.5	8168.0	8168.5	7641.0	7650.5	7666.5	7666.5	8168.0	8168.5	7614.5	8095.5	8102.0	12.0	
(0.75, 0.20)	13928.0	11446.0	11446.0	5947.0	6036.5	6245.0	2106.0	8675.0	8753.0	8774.0	8957.0	10772.0	11149.0	8675.0	8753.0	8847.0	8994.0	10772.0	11149.0	8137.0	9813.5	9947.5	11.0	

Table 24: Median number of iterations: Model 4 (high-dimensional, slowly decaying coefficients)

1.3 Median number of predictors

The results shown in this section are similar to those presented in Section 1.2, except that this time we report the median number of predictors selected when various stopping rules are applied. The conventions for the type and the colour of the font used in Tables 25-36 are the same as in Tables 13-24.

1.3.1 Experiments for $n = 20$

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	30.0	42.0	29.0	42.0	42.0	42.0	41.0	34.0	34.0	33.0	33.0	8.0	7.0	34.0	34.0	33.0	33.0	9.0	9.0	34.0	40.0	40.0	25.0
(0.00, 0.20)	30.0	40.5	28.0	40.5	40.5	40.5	40.0	36.0	36.0	35.0	35.0	1.0	1.0	35.5	35.5	34.0	34.5	1.0	1.0	35.5	40.0	40.0	4.0
(0.75, 8.00)	27.5	43.0	24.0	43.0	43.0	43.0	40.0	36.0	36.0	30.5	31.0	4.0	4.0	35.0	36.0	30.0	31.0	5.0	5.0	36.0	42.0	42.0	9.0
(0.75, 0.20)	26.5	41.0	22.0	41.0	41.0	41.0	37.0	39.0	39.0	39.0	39.0	1.0	1.0	39.0	39.0	39.0	39.0	1.0	1.0	39.0	41.0	41.0	3.0
Case 2																							
(0.00, 8.00)	31.0	44.0	26.5	44.0	44.0	44.0	41.0	37.0	37.0	35.0	35.5	7.5	7.5	37.0	37.0	35.5	35.5	8.0	8.0	37.0	44.0	44.0	5.0
(0.00, 0.20)	31.5	43.5	27.5	43.5	43.5	43.5	42.0	40.0	40.0	38.5	38.5	6.0	6.0	39.0	39.0	38.0	38.0	6.0	6.0	39.0	43.0	43.0	5.0
(0.75, 8.00)	29.0	45.5	21.0	45.5	45.5	45.5	9.0	15.5	15.5	13.0	13.0	7.0	7.0	16.5	16.5	13.0	13.0	7.0	7.5	16.5	43.0	43.0	4.0
(0.75, 0.20)	28.0	45.0	18.5	45.0	45.0	45.0	6.5	41.0	41.0	39.0	39.0	6.0	6.0	41.0	41.0	39.0	39.0	6.0	6.0	41.0	45.0	45.0	5.0
Case 3																							
(0.00, 8.00)	31.0	45.0	27.0	45.0	45.0	45.0	43.0	37.0	38.0	34.0	34.5	8.0	8.0	36.0	36.0	34.5	34.5	9.0	9.0	37.0	44.0	44.0	12.0
(0.00, 0.20)	31.0	44.0	26.0	44.0	44.0	44.0	40.0	39.0	39.0	38.5	38.5	5.0	5.0	39.0	39.0	38.0	38.0	6.0	6.0	39.0	43.5	43.5	9.0
(0.75, 8.00)	27.0	45.0	20.5	45.0	45.0	45.0	21.0	26.0	26.0	14.0	14.0	6.0	6.0	26.0	26.0	14.0	16.0	6.0	6.5	26.0	45.0	45.0	7.5
(0.75, 0.20)	27.0	45.0	20.0	45.0	45.0	45.0	6.5	40.0	40.0	39.0	39.5	3.5	3.0	40.0	40.0	39.0	39.0	4.0	4.0	40.0	44.5	44.5	10.0

Table 25: Median number of predictors: Model 1 (low-dimensional)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	31.0	42.0	30.0	42.0	42.0	42.0	40.5	36.0	36.5	35.0	35.0	1.0	1.0	36.0	36.0	35.0	35.0	1.0	2.0	36.0	41.5	41.5	17.0
(0.00, 0.20)	29.0	41.0	28.0	41.0	41.0	41.0	39.0	36.0	36.0	35.0	35.0	1.0	1.0	35.5	35.5	35.0	35.0	1.0	1.0	35.5	40.0	40.0	6.0
(0.75, 8.00)	31.0	47.0	27.0	47.0	47.0	47.0	42.0	36.0	36.5	34.0	35.0	8.0	7.0	36.0	36.0	34.0	34.5	10.0	10.0	36.0	47.0	47.0	42.0
(0.75, 0.20)	26.0	42.0	21.5	42.0	42.0	42.0	34.5	39.0	39.0	39.0	39.0	1.0	1.0	39.0	39.0	38.5	38.5	1.0	1.0	39.0	41.0	41.0	4.5
Case 2																							
(0.00, 8.00)	30.0	44.0	24.0	44.0	44.0	44.0	41.5	39.0	39.0	38.0	38.0	7.0	7.0	39.0	39.0	38.0	38.0	7.0	7.0	39.0	44.0	44.0	5.0
(0.00, 0.20)	31.0	43.0	24.0	43.0	43.0	43.0	39.0	39.0	39.0	38.0	38.0	6.0	6.0	39.0	39.0	38.0	38.0	7.0	6.0	39.0	43.0	43.0	5.0
(0.75, 8.00)	31.0	47.0	25.5	47.0	47.0	47.0	10.0	33.5	33.5	23.0	25.0	8.0	8.0	33.5	33.5	21.5	23.5	8.5	8.5	33.5	47.0	47.0	7.0
(0.75, 0.20)	28.0	45.0	18.0	45.0	45.0	45.0	6.0	40.0	40.0	37.0	37.0	6.0	6.0	39.5	40.0	37.0	37.0	6.0	6.0	40.0	45.0	45.0	5.5
Case 3																							
(0.00, 8.00)	31.5	44.0	26.0	44.0	44.0	44.0	42.0	39.0	39.0	37.5	38.0	6.0	6.0	39.0	39.0	37.0	37.0	6.0	6.0	39.0	44.0	44.0	7.5
(0.00, 0.20)	30.0	43.0	27.0	43.0	43.0	43.0	40.0	38.0	38.0	38.0	38.0	5.0	5.0	38.0	38.0	38.0	38.0	5.5	5.5	38.0	42.0	42.0	6.0
(0.75, 8.00)	32.0	49.0	27.5	49.0	49.0	49.0	43.0	39.5	40.0	36.0	37.0	8.0	8.0	39.5	39.5	35.5	36.5	8.0	8.0	39.5	48.0	48.0	18.0
(0.75, 0.20)	26.0	45.0	20.5	45.0	45.0	45.0	2.5	41.0	41.0	39.5	39.5	4.0	3.0	40.5	40.5	39.0	39.0	4.0	4.0	40.5	45.0	45.0	4.5

Table 26: Median number of predictors: Model 2 (high-dimensional, small equal coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	31.0	41.0	29.0	41.0	41.0	41.0	40.0	36.0	36.0	33.0	34.5	4.5	4.0	35.0	35.0	33.0	33.0	6.0	5.0	35.5	40.0	40.0	14.0
(0.00, 0.20)	30.0	39.5	28.0	39.5	39.5	39.5	39.0	35.0	35.0	34.0	34.0	1.0	1.0	35.0	35.0	34.0	34.0	1.0	1.0	35.0	39.0	39.0	4.5
(0.75, 8.00)	30.0	46.0	26.0	46.0	46.0	46.0	42.0	36.0	36.0	32.5	33.5	8.0	8.0	35.0	36.0	32.0	33.0	8.0	8.0	36.0	45.0	45.0	24.5
(0.75, 0.20)	26.0	41.0	23.0	41.0	41.0	41.0	32.0	38.0	38.0	37.5	37.5	1.0	1.0	38.0	38.0	37.0	37.0	1.0	1.0	38.0	40.5	40.5	4.5
Case 2																							
(0.00, 8.00)	30.5	44.0	25.0	44.0	44.0	44.0	40.5	34.0	34.5	28.5	29.0	7.0	7.0	33.0	33.5	28.5	29.0	7.0	7.0	34.0	44.0	44.0	4.0
(0.00, 0.20)	31.0	43.0	25.0	43.0	43.0	43.0	38.0	38.0	38.0	37.0	37.5	6.0	6.0	38.0	38.0	36.0	36.5	7.0	7.0	38.0	43.0	43.0	4.0
(0.75, 8.00)	30.0	47.0	24.0	47.0	47.0	47.0	41.0	22.0	23.5	15.5	15.5	8.0	8.0	22.0	23.5	15.5	15.5	8.0	8.0	23.5	45.5	46.0	4.0
(0.75, 0.20)	27.0	46.0	19.5	46.0	46.0	46.0	6.0	40.5	41.0	38.5	38.5	6.0	6.0	40.0	40.0	38.5	38.5	6.0	6.0	40.5	45.0	45.0	6.0
Case 3																							
(0.00, 8.00)	31.0	45.0	25.5	45.0	45.0	45.0	42.0	39.0	39.0	37.5	38.0	6.0	6.0	39.0	39.0	37.0	37.5	7.0	7.0	39.0	45.0	45.0	9.0
(0.00, 0.20)	29.5	44.0	24.5	44.0	44.0	44.0	39.0	40.0	40.0	39.0	39.0	4.5	4.0	40.0	40.0	39.0	39.0	5.0	5.0	40.0	43.0	43.0	6.5
(0.75, 8.00)	30.0	47.0	23.0	47.0	47.0	47.0	39.5	29.0	29.0	19.0	20.0	9.0	9.0	29.0	29.0	19.5	20.0	9.0	9.0	29.0	47.0	47.0	12.5
(0.75, 0.20)	27.0	46.0	20.0	46.0	46.0	46.0	4.5	41.0	41.5	40.0	40.0	4.5	4.0	41.0	41.0	39.5	40.0	5.0	5.0	41.0	45.0	45.0	8.5

Table 27: Median number of predictors: Model 3 (high-dimensional, decaying coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	30.0	41.5	27.5	41.5	41.5	41.5	40.0	38.0	38.0	36.0	36.5	1.0	1.0	37.5	38.0	36.0	36.0	1.0	2.0	38.0	41.0	41.0	10.5
(0.00, 0.20)	30.0	40.5	28.0	40.5	40.5	40.5	39.0	37.0	37.0	36.0	36.0	1.0	1.0	36.0	36.0	36.0	36.0	1.0	1.0	36.0	40.0	40.0	3.0
(0.75, 8.00)	32.0	48.0	29.0	48.0	48.0	48.0	45.0	40.0	40.0	35.0	35.5	8.0	8.5	39.0	39.5	34.0	35.0	9.5	10.0	39.5	47.5	47.5	41.0
(0.75, 0.20)	26.0	42.0	21.0	42.0	42.0	42.0	1.0	39.0	39.0	39.0	39.0	1.0	1.0	39.0	39.0	38.5	39.0	1.0	1.0	39.0	42.0	42.0	5.0
Case 2																							
(0.00, 8.00)	30.0	44.0	26.0	44.0	44.0	44.0	41.0	39.0	39.0	38.0	38.0	6.0	6.0	39.0	39.0	37.5	38.0	6.0	6.0	39.0	43.0	43.0	7.5
(0.00, 0.20)	30.5	44.0	26.0	44.0	44.0	44.0	40.0	38.0	38.0	37.0	37.0	7.0	7.0	38.0	38.0	37.0	37.0	7.0	7.0	38.0	44.0	44.0	3.0
(0.75, 8.00)	31.5	48.0	26.0	48.0	48.0	48.0	37.0	34.5	35.0	22.5	23.5	8.0	8.0	34.0	34.0	22.5	23.5	8.5	9.0	34.0	47.0	47.0	7.0
(0.75, 0.20)	27.0	45.0	19.0	45.0	45.0	45.0	7.5	38.0	38.0	37.0	37.0	6.0	6.0	38.0	38.0	37.0	37.0	6.0	6.0	38.0	44.0	44.0	7.0
Case 3																							
(0.00, 8.00)	31.0	45.0	28.0	45.0	45.0	45.0	41.5	39.0	39.5	37.5	37.5	6.0	5.0	39.0	39.0	37.0	37.0	6.0	6.0	39.0	45.0	45.0	11.0
(0.00, 0.20)	30.0	43.0	26.0	43.0	43.0	43.0	40.0	39.0	39.0	39.0	39.0	5.5	5.0	39.0	39.0	39.0	39.0	6.0	6.0	39.0	43.0	43.0	8.5
(0.75, 8.00)	31.0	49.0	27.0	49.0	49.0	49.0	44.0	40.0	40.0	31.5	34.5	8.0	8.0	38.5	40.0	31.5	33.5	8.0	8.0	40.0	49.0	49.0	19.5
(0.75, 0.20)	26.0	46.0	19.0	46.0	46.0	46.0	4.0	42.0	42.0	41.0	41.0	4.0	3.0	42.0	42.0	41.0	41.0	4.0	4.0	42.0	45.0	45.0	8.0

Table 28: Median number of predictors: Model 4 (high-dimensional, slowly decaying coefficients)

1.3.2 Experiments for $n = 100$

(ω, ς^2)	AIC_C	KIC	KIC_C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	45.5	98.0	32.5	27.5	22.5	8.0	4.0	21.0	20.0	15.0	14.0	7.0	7.0	22.0	21.0	16.0	15.0	7.0	7.0	21.0	28.0	27.0	12.0
(0.00, 0.20)	45.0	98.0	31.5	23.0	16.5	2.0	1.0	53.0	50.0	54.5	52.5	1.0	1.0	53.0	51.5	55.0	53.5	3.0	2.0	51.0	98.0	98.0	8.5
(0.75, 8.00)	40.0	92.0	21.0	14.5	11.0	5.5	4.0	10.0	10.0	7.0	7.0	4.0	4.0	11.0	11.0	8.0	8.0	5.0	5.0	11.0	10.0	10.0	8.0
(0.75, 0.20)	37.0	91.0	20.0	12.5	8.0	3.0	1.0	63.5	55.0	76.0	76.0	3.0	2.0	65.0	56.0	77.5	78.0	3.0	3.0	62.5	93.0	93.0	6.0
Case 2																							
(0.00, 8.00)	68.0	85.0	60.0	66.0	70.5	40.0	16.5	42.0	42.0	41.0	41.0	28.0	28.0	42.0	42.0	41.0	41.0	28.0	28.0	42.0	47.5	47.5	3.0
(0.00, 0.20)	65.0	84.0	55.0	56.5	59.0	22.0	11.0	46.0	46.0	43.5	45.0	18.5	19.0	46.0	46.0	45.5	45.5	19.0	20.0	46.0	57.0	57.5	3.0
(0.75, 8.00)	68.0	77.0	57.0	61.0	66.0	24.0	11.0	33.5	33.5	30.0	30.0	17.0	17.0	34.0	34.0	31.0	31.0	17.0	17.0	34.0	37.0	37.0	2.0
(0.75, 0.20)	66.0	76.5	54.5	56.0	63.0	20.0	9.0	44.0	44.0	42.5	42.5	18.0	17.5	45.0	45.0	43.0	43.0	18.0	18.0	45.0	62.0	62.0	2.0
Case 3																							
(0.00, 8.00)	58.0	92.0	47.0	47.0	47.0	23.0	13.0	32.5	33.0	29.0	29.0	20.0	20.0	33.0	33.0	30.0	30.0	20.0	20.0	33.0	38.0	38.0	6.0
(0.00, 0.20)	56.0	91.0	45.0	42.0	38.0	17.0	7.5	46.0	46.0	45.0	44.5	17.0	17.0	47.0	46.5	45.0	46.0	17.5	17.0	46.0	66.0	66.0	11.0
(0.75, 8.00)	62.0	85.0	44.5	41.0	38.5	16.0	9.0	23.5	23.5	20.0	20.5	14.0	14.0	24.0	24.0	22.0	22.0	14.0	14.0	24.0	28.0	28.0	3.0
(0.75, 0.20)	61.5	85.0	44.0	37.5	33.5	13.0	7.0	49.0	49.0	47.0	45.5	13.0	13.0	50.0	49.5	47.5	47.0	13.5	13.0	49.0	75.5	76.0	8.0

Table 29: Median number of predictors: Model 1 (low-dimensional)

(ω, ς^2)	AIC_C	KIC	KIC_C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	70.0	99.0	60.0	99.0	99.0	99.0	1.0	69.0	71.0	73.0	77.0	19.0	10.0	69.0	71.0	73.0	77.0	22.0	12.0	73.0	85.0	88.0	88.0
(0.00, 0.20)	45.0	98.0	33.0	26.0	19.0	1.0	1.0	53.0	53.0	55.0	55.0	1.0	1.0	54.0	53.0	55.5	55.5	1.0	1.0	53.0	97.0	98.0	11.0
(0.75, 8.00)	64.0	95.0	57.0	57.0	57.5	50.0	14.0	53.0	54.0	54.0	55.0	47.0	48.0	53.5	54.5	54.0	55.0	47.0	48.0	54.5	57.0	58.5	69.0
(0.75, 0.20)	36.5	92.0	23.0	18.0	14.5	4.0	1.0	60.0	51.0	74.5	70.0	5.0	2.0	61.0	53.5	75.0	70.0	6.0	5.0	48.5	93.0	93.5	12.5
Case 2																							
(0.00, 8.00)	83.0	89.0	74.0	89.0	89.0	88.0	12.0	64.0	64.0	67.0	67.0	26.0	26.0	64.0	64.0	67.5	67.5	26.5	26.5	64.0	75.0	75.0	3.0
(0.00, 0.20)	67.0	85.0	56.0	60.0	65.5	22.5	10.5	49.0	49.0	49.0	49.0	20.0	19.5	49.0	49.0	49.0	49.0	20.0	20.0	49.0	62.0	61.0	3.0
(0.75, 8.00)	76.0	81.0	70.0	77.0	79.0	63.0	15.0	57.5	58.0	59.0	59.0	43.0	43.0	58.0	58.0	59.0	59.0	43.0	43.5	58.0	64.0	64.0	3.0
(0.75, 0.20)	68.0	76.5	54.0	59.5	64.5	20.0	10.0	48.0	48.0	47.0	46.5	16.0	16.0	48.0	48.5	47.5	47.0	16.5	16.5	48.5	64.0	64.0	3.0
Case 3																							
(0.00, 8.00)	80.0	96.0	70.0	95.0	96.0	95.0	9.0	68.0	69.5	72.5	73.0	27.0	27.0	68.5	70.0	73.0	73.0	27.0	27.0	70.0	78.5	79.5	30.5
(0.00, 0.20)	58.0	91.0	47.5	45.5	45.0	16.5	8.0	47.0	47.0	47.0	46.5	16.0	16.0	47.5	47.0	48.0	47.0	17.0	17.0	47.0	65.0	65.0	12.5
(0.75, 8.00)	73.0	89.0	64.0	71.0	77.5	52.5	18.0	56.5	57.0	58.0	58.0	44.0	44.5	57.0	57.0	58.0	58.0	44.0	44.5	57.0	61.5	61.5	29.5
(0.75, 0.20)	58.0	85.0	43.5	38.0	33.5	13.0	7.0	46.5	45.0	45.5	43.5	13.0	13.0	47.5	46.0	45.5	45.5	13.0	13.0	46.0	74.5	75.0	12.5

Table 30: Median number of predictors: Model 2 (high-dimensional, small equal coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ ⁺	ESC ₂ ⁺	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ ⁺	EgMDL ₂ ⁺	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	54.5	98.5	44.0	45.0	43.5	17.0	6.0	39.0	38.5	37.0	37.0	15.0	14.0	39.0	38.5	37.0	37.0	15.0	15.0	38.5	49.0	49.0	35.5
(0.00, 0.20)	44.0	98.0	31.0	21.5	14.0	2.0	1.0	54.0	49.0	56.5	52.5	1.0	1.0	55.0	49.5	57.0	54.0	3.0	2.0	49.0	96.0	98.0	7.5
(0.75, 8.00)	47.0	92.0	37.0	33.0	32.0	22.0	13.0	30.0	30.5	28.5	28.5	21.0	21.0	30.0	30.5	29.0	29.0	21.0	21.0	30.5	33.0	33.0	35.0
(0.75, 0.20)	39.0	91.5	22.0	15.0	10.0	3.0	1.0	60.0	56.5	75.0	74.5	3.0	2.0	60.0	58.0	76.0	76.0	4.0	4.0	57.5	93.0	93.0	8.0
Case 2																							
(0.00, 8.00)	71.5	86.0	65.0	77.5	82.5	53.5	14.0	50.0	50.0	50.0	50.0	31.0	31.0	50.0	50.0	50.0	50.0	31.5	31.5	50.0	55.0	55.0	3.0
(0.00, 0.20)	65.0	85.0	54.0	58.0	62.0	21.0	11.0	47.0	47.0	46.0	46.0	18.0	18.5	47.0	47.0	47.0	46.5	19.0	19.0	47.0	57.5	57.5	3.0
(0.75, 8.00)	70.0	78.0	60.0	70.0	73.0	40.0	14.0	42.5	42.5	42.0	42.0	25.5	25.5	43.0	43.0	42.0	42.0	25.5	26.0	43.0	48.0	48.0	2.0
(0.75, 0.20)	67.0	77.0	56.0	61.5	65.0	20.5	9.0	49.0	50.0	49.0	50.0	17.5	17.5	50.0	50.0	50.0	50.0	18.5	18.0	50.0	63.5	63.0	3.0
Case 3																							
(0.00, 8.00)	65.5	94.0	58.0	67.5	76.0	34.0	13.0	46.0	46.0	44.0	44.0	25.0	24.5	46.0	46.0	44.0	44.5	26.0	26.0	46.0	55.0	55.0	10.5
(0.00, 0.20)	56.5	91.5	46.0	41.5	41.0	15.0	8.0	46.5	46.0	46.0	45.0	15.0	15.0	47.0	47.5	46.0	46.0	16.0	15.0	46.0	66.0	66.0	13.0
(0.75, 8.00)	62.5	86.5	51.0	52.0	52.0	25.0	15.0	33.5	33.5	32.0	32.0	21.0	21.5	34.0	34.0	32.0	32.0	22.0	22.0	34.0	41.5	41.5	6.5
(0.75, 0.20)	58.5	86.0	43.0	40.5	37.0	13.0	7.0	44.0	44.0	41.0	41.0	13.0	13.0	44.5	45.0	41.0	41.0	13.0	13.0	43.5	73.0	73.0	9.5

Table 31: Median number of predictors: Model 3 (high-dimensional, decaying coefficients)

(ω, ς^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ ⁺	ESC ₂ ⁺	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ ⁺	EgMDL ₂ ⁺	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	68.0	99.0	59.0	95.0	98.5	99.0	2.0	63.0	65.5	66.0	69.0	23.0	19.5	64.0	66.0	67.0	69.0	24.0	21.0	66.0	76.0	78.0	84.5
(0.00, 0.20)	45.0	98.0	32.5	23.0	18.0	1.0	1.0	56.5	52.0	59.0	57.0	1.0	1.0	57.0	55.0	60.0	59.0	1.0	2.0	55.0	97.5	98.0	7.5
(0.75, 8.00)	58.0	94.0	53.0	51.0	51.0	44.0	16.0	48.5	49.0	48.0	49.0	41.0	42.0	49.0	49.0	49.0	49.0	41.5	42.0	49.0	51.5	52.0	59.0
(0.75, 0.20)	39.0	91.0	23.0	16.0	13.0	3.0	1.0	58.0	50.5	72.0	62.0	3.5	1.5	58.5	52.5	72.0	72.0	5.0	4.0	49.0	93.0	93.0	11.0
Case 2																							
(0.00, 8.00)	79.5	87.0	74.0	87.0	87.0	86.0	12.0	61.0	62.0	66.0	65.5	31.0	31.0	61.5	62.0	66.0	66.0	31.5	31.5	62.0	69.5	70.0	3.0
(0.00, 0.20)	66.5	85.0	56.5	61.5	65.5	22.0	11.0	48.5	48.5	48.5	48.0	19.0	19.0	50.0	50.0	50.0	50.0	19.0	19.0	50.0	61.0	61.0	2.5
(0.75, 8.00)	75.0	80.0	67.5	75.0	76.0	53.0	17.0	53.0	53.0	53.0	53.0	39.5	39.5	53.0	53.0	53.0	53.0	39.5	39.5	53.0	58.0	58.0	3.0
(0.75, 0.20)	68.0	77.0	57.0	59.0	64.5	18.0	9.0	47.0	47.0	46.5	46.0	15.5	15.5	48.0	48.0	47.5	47.0	16.0	16.0	47.0	64.0	64.0	2.0
Case 3																							
(0.00, 8.00)	75.0	95.0	66.0	94.0	95.0	94.0	10.0	61.0	61.0	64.0	64.5	29.0	29.0	61.0	62.0	64.5	64.5	30.5	30.0	62.0	73.0	73.0	25.5
(0.00, 0.20)	58.5	93.0	46.0	43.5	44.0	15.0	7.5	47.0	47.0	46.0	45.5	15.0	15.0	47.5	47.5	46.5	46.5	15.0	15.0	48.0	70.0	69.0	11.0
(0.75, 8.00)	69.0	88.0	61.0	65.0	69.0	42.0	17.0	49.0	49.5	49.0	49.0	35.5	36.0	49.5	49.5	49.5	49.0	36.0	36.0	49.5	55.0	56.0	17.0
(0.75, 0.20)	55.5	84.0	41.0	35.0	31.0	14.0	7.0	42.5	42.0	37.5	37.0	14.0	14.0	44.5	42.5	38.5	38.5	14.0	14.0	42.0	75.5	75.5	11.5

Table 32: Median number of predictors: Model 4 (high-dimensional, slowly decaying coefficients)

1.3.3 Experiments for $n = 10000$

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	42.0	31.0	31.0	10.0	7.0	6.0	4.0	9.0	9.0	7.0	7.0	5.0	5.0	9.0	9.0	7.0	7.0	6.0	6.0	9.0	9.0	9.0	46.0
(0.00, 0.20)	42.0	31.0	31.0	9.0	7.0	5.0	4.0	15.0	15.0	11.0	10.0	6.0	6.0	16.0	16.0	11.0	11.0	7.0	7.0	11.0	17.0	16.0	97.0
(0.75, 8.00)	35.0	20.0	20.0	7.0	5.0	4.0	3.0	6.0	6.0	5.0	5.0	4.0	4.0	6.0	6.0	5.0	5.0	4.0	4.0	6.0	5.0	5.0	26.0
(0.75, 0.20)	35.5	20.0	20.0	7.0	5.0	4.0	3.0	9.0	9.0	6.0	6.0	4.0	4.0	9.0	9.0	7.0	7.0	4.0	4.0	7.0	7.0	7.0	89.0
Case 2																							
(0.00, 8.00)	92.5	90.0	90.0	77.0	77.0	76.0	52.0	79.0	79.0	80.0	80.0	79.5	80.0	79.0	79.0	80.0	80.0	80.0	80.0	79.0	81.0	81.0	3.0
(0.00, 0.20)	92.0	90.0	90.0	76.5	77.0	76.0	51.5	86.0	86.0	87.5	88.0	91.0	100.0	86.0	86.0	87.5	88.0	91.0	100.0	84.0	88.0	89.0	3.0
(0.75, 8.00)	96.0	94.0	94.0	83.0	84.0	83.0	51.0	86.0	86.0	87.0	87.0	87.0	87.0	86.0	86.0	87.0	87.0	87.0	87.0	86.0	87.0	87.5	2.0
(0.75, 0.20)	96.0	94.0	94.0	82.5	83.5	83.0	50.0	91.0	91.0	93.0	93.0	95.0	97.0	91.0	91.0	93.0	93.0	95.0	97.0	89.0	93.0	93.0	2.0
Case 3																							
(0.00, 8.00)	94.0	92.0	92.0	82.0	83.0	82.0	65.0	83.0	83.0	84.0	84.0	83.5	84.0	83.0	83.0	84.0	84.0	84.0	84.0	83.0	85.0	85.0	3.0
(0.00, 0.20)	93.5	92.0	92.0	82.0	82.5	82.0	65.0	87.0	87.0	88.5	89.0	89.0	89.0	87.0	87.0	89.0	89.0	89.0	89.0	86.0	89.0	89.0	7.0
(0.75, 8.00)	97.0	96.0	96.0	88.0	89.0	89.0	70.0	90.0	90.0	91.0	91.0	91.0	91.0	90.0	90.0	91.0	91.0	91.0	91.0	90.0	91.0	91.0	2.0
(0.75, 0.20)	97.0	95.0	95.0	88.0	89.0	88.5	69.0	93.0	93.0	94.0	94.0	94.5	95.0	93.0	93.0	94.0	94.0	95.0	95.0	92.0	94.0	94.0	3.0

Table 33: Median number of predictors: Model 1 (low-dimensional)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	99.0
(0.00, 0.20)	100.0	100.0	100.0	99.0	100.0	99.0	81.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	99.0
(0.75, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	99.0
(0.75, 0.20)	79.0	78.0	78.0	75.0	76.0	74.0	68.0	77.0	77.5	78.0	78.0	76.0	76.0	77.0	77.0	78.0	78.0	78.0	76.0	76.0	77.0	78.0	98.0
Case 2																							
(0.00, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	50.0
(0.00, 0.20)	96.0	94.5	94.0	87.0	88.0	88.0	70.0	92.0	93.0	93.0	94.0	100.0	100.0	92.0	93.0	93.0	94.0	100.0	100.0	91.0	94.0	94.0	23.5
(0.75, 8.00)	99.0	98.0	98.0	95.0	95.0	95.0	91.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	96.0	25.0
(0.75, 0.20)	95.0	93.0	93.0	81.5	82.0	81.0	54.0	90.0	90.0	92.0	92.0	93.5	94.0	90.0	90.0	92.0	92.0	94.0	94.0	88.0	92.0	93.0	13.0
Case 3																							
(0.00, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	10.5
(0.00, 0.20)	96.0	95.0	95.0	88.0	89.0	89.0	73.0	93.0	93.0	93.0	94.0	95.0	96.0	93.0	93.0	93.0	94.0	95.0	96.0	92.0	94.0	94.0	5.5
(0.75, 8.00)	99.0	98.5	98.5	96.0	96.0	96.0	91.0	96.0	96.0	97.0	97.0	97.0	97.0	96.0	96.0	97.0	97.0	97.0	97.0	96.0	97.0	97.0	15.0
(0.75, 0.20)	97.0	95.0	95.0	88.0	89.0	88.0	70.0	92.0	93.0	94.0	94.0	95.0	95.0	92.0	93.0	94.0	94.0	95.0	95.0	92.0	94.0	94.0	5.5

Table 34: Median number of predictors: Model 2 (high-dimensional, small equal coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	99.0	99.0	99.0	97.0	98.0	98.0	78.0	98.0	98.0	99.0	99.0	99.0	99.0	98.0	98.0	99.0	99.0	99.0	99.0	98.0	99.0	99.0	99.0
(0.00, 0.20)	65.5	57.0	57.0	31.0	29.0	23.0	13.0	47.0	47.0	47.0	46.0	32.5	32.0	47.0	47.0	47.0	46.5	33.0	32.5	38.0	53.0	53.0	97.0
(0.75, 8.00)	87.0	87.0	87.0	86.0	86.0	85.0	82.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	86.0	93.0
(0.75, 0.20)	50.0	45.0	44.5	35.0	34.5	31.0	23.0	39.0	39.5	39.0	39.0	34.0	34.0	40.0	40.0	39.0	39.0	34.0	34.0	37.5	41.0	41.0	49.0
Case 2																							
(0.00, 8.00)	97.0	96.0	96.0	90.0	91.0	92.0	79.0	92.0	92.0	92.0	92.5	94.0	94.0	92.0	92.0	92.0	92.5	94.0	94.0	92.0	93.0	93.0	1.0
(0.00, 0.20)	93.0	91.0	91.0	78.0	79.0	77.5	54.0	87.5	88.0	89.0	89.0	100.0	100.0	88.0	88.0	89.0	89.0	100.0	100.0	85.0	90.0	90.0	1.0
(0.75, 8.00)	96.0	94.0	94.0	84.0	85.0	83.5	64.0	86.0	86.0	87.5	87.5	86.0	86.0	86.0	86.0	87.5	87.5	86.0	86.0	86.0	87.0	87.0	3.0
(0.75, 0.20)	95.0	93.0	93.0	81.5	82.5	82.0	52.5	90.0	90.0	91.5	92.0	94.0	95.0	90.0	90.0	92.0	92.0	94.0	95.0	88.0	92.0	93.0	2.5
Case 3																							
(0.00, 8.00)	98.0	97.0	97.0	93.0	93.0	93.0	83.0	93.0	93.0	94.0	94.0	95.0	95.0	93.0	93.0	94.0	94.0	95.0	95.0	93.0	94.0	94.0	1.0
(0.00, 0.20)	94.0	92.0	92.0	84.0	84.0	84.0	68.0	88.0	88.0	89.0	89.0	90.0	90.0	88.0	88.5	89.0	89.0	90.0	90.0	88.0	90.0	90.0	5.0
(0.75, 8.00)	97.0	96.0	96.0	90.0	91.0	91.0	76.0	91.0	91.0	92.0	92.0	92.0	92.0	91.0	91.0	92.0	92.0	92.0	92.0	91.0	92.0	92.0	3.0
(0.75, 0.20)	97.0	96.0	96.0	88.0	89.5	89.0	71.0	93.0	93.0	94.0	94.0	95.0	95.0	93.0	93.0	94.0	94.0	95.0	95.0	92.0	94.0	94.0	4.0

Table 35: Median number of predictors: Model 3 (high-dimensional, decaying coefficients)

(ω, ζ^2)	AIC _C	KIC	KIC _C	BIC	EBIC	EBIC*	EBIC ^o	SC ₁	SC ₂	ESC ₁	ESC ₂	ESC ₁ *	ESC ₂ *	gMDL ₁	gMDL ₂	EgMDL ₁	EgMDL ₂	EgMDL ₁ *	EgMDL ₂ *	MMLU	MMLG ₁	MMLG ₂	CV
Case 1																							
(0.00, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	99.0
(0.00, 0.20)	98.0	97.0	97.0	85.5	86.0	83.0	42.0	96.0	97.0	97.0	98.0	100.0	100.0	96.0	97.0	97.0	98.0	100.0	100.0	94.0	98.0	98.0	99.0
(0.75, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	99.0
(0.75, 0.20)	71.0	69.0	69.0	66.0	66.0	64.0	58.0	68.0	68.0	69.0	69.0	67.0	67.0	68.0	69.0	69.0	69.0	67.0	67.0	67.5	69.0	69.0	98.0
Case 2																							
(0.00, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	99.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	6.0
(0.00, 0.20)	95.0	94.0	94.0	84.0	85.0	85.0	65.0	91.0	91.0	92.0	93.0	100.0	100.0	91.0	91.0	92.0	93.0	100.0	100.0	89.0	93.0	93.0	11.5
(0.75, 8.00)	97.0	96.0	96.0	90.0	91.0	90.0	83.0	92.0	92.0	93.0	93.0	91.0	91.0	92.0	92.0	93.0	93.0	91.0	91.0	92.0	93.0	93.0	7.0
(0.75, 0.20)	95.0	93.0	93.0	82.0	83.0	82.0	54.0	91.0	91.0	92.0	92.0	94.0	95.0	91.0	91.0	92.0	92.0	94.0	95.0	89.0	92.0	93.0	8.5
Case 3																							
(0.00, 8.00)	100.0	100.0	100.0	100.0	100.0	100.0	99.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	7.0
(0.00, 0.20)	96.0	94.0	94.0	87.0	88.0	87.5	72.0	91.0	91.0	92.0	92.0	93.0	93.0	91.0	91.0	92.0	92.0	93.0	93.0	90.5	92.0	92.0	5.0
(0.75, 8.00)	99.0	98.0	98.0	94.0	94.5	94.0	86.0	95.0	95.0	95.0	95.0	95.0	95.0	95.0	95.0	95.0	95.0	95.0	95.0	94.0	95.0	95.0	6.5
(0.75, 0.20)	97.0	95.0	95.0	88.0	89.0	89.0	71.5	92.0	92.0	93.5	94.0	94.5	95.0	92.0	92.0	93.5	94.0	94.5	95.0	92.0	93.5	93.5	5.0

Table 36: Median number of predictors: Model 4 (high-dimensional, slowly decaying coefficients)

1.4 Score plots

In this section, we plot graphs similar to those in [1, Figs. 1-4]. The scores presented in each graph are computed by aggregating the scores obtained from all the experiments that have a common parameter. For example, we show in Figure 1 the scores aggregated from all the experiments in which SNR is high ($\zeta^2 = 8$).

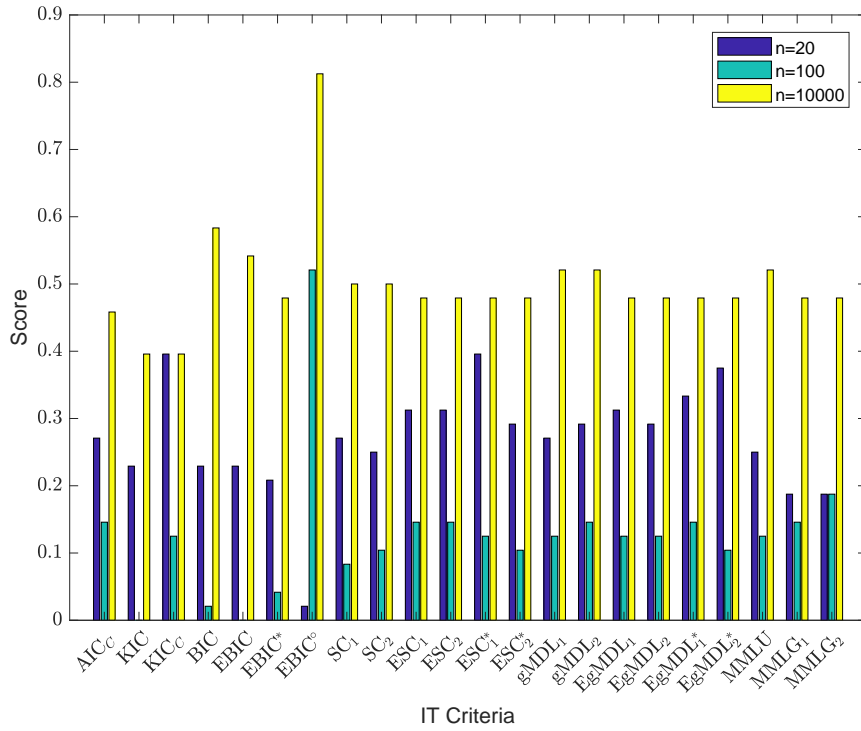


Figure 1: Scores aggregated from all experiments in which SNR is high ($\zeta^2 = 8$).

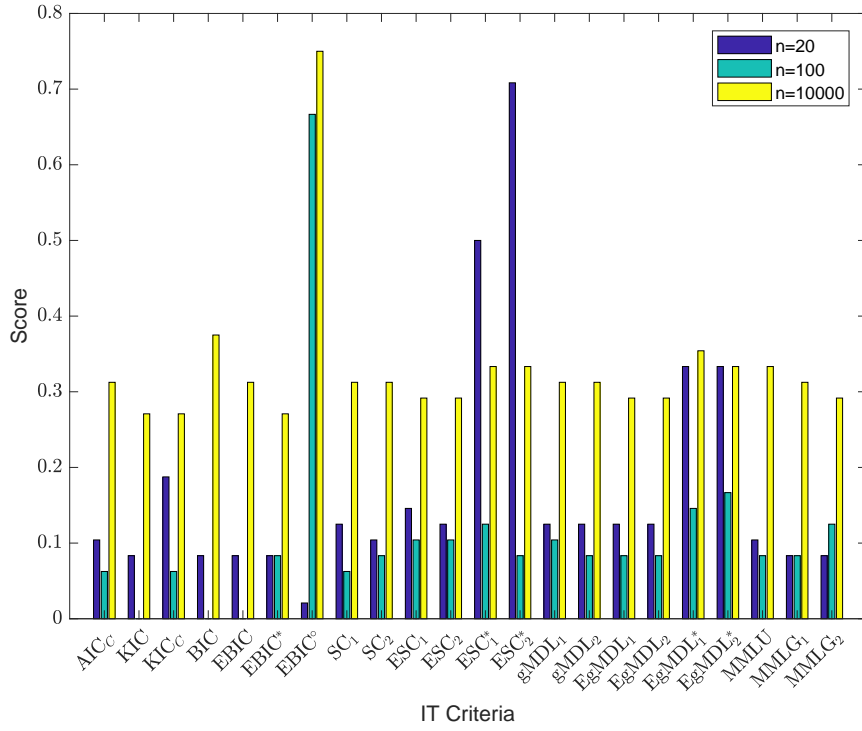


Figure 2: Scores aggregated from all experiments in which $\omega = 0$.

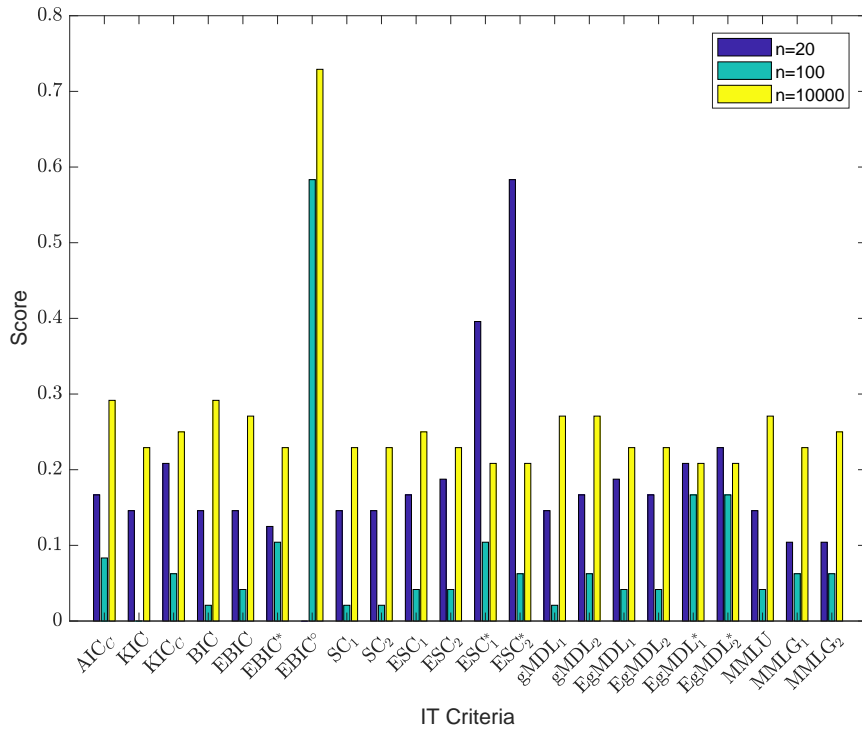


Figure 3: Scores aggregated from all experiments in which $\omega = 0.75$.

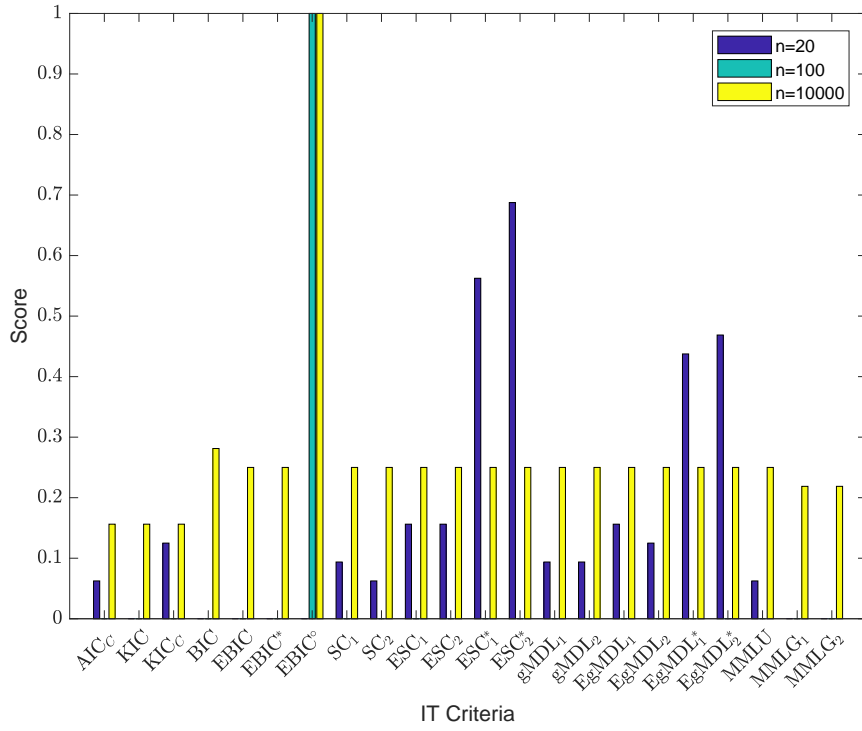


Figure 4: Scores aggregated from all experiments done for Case 2.

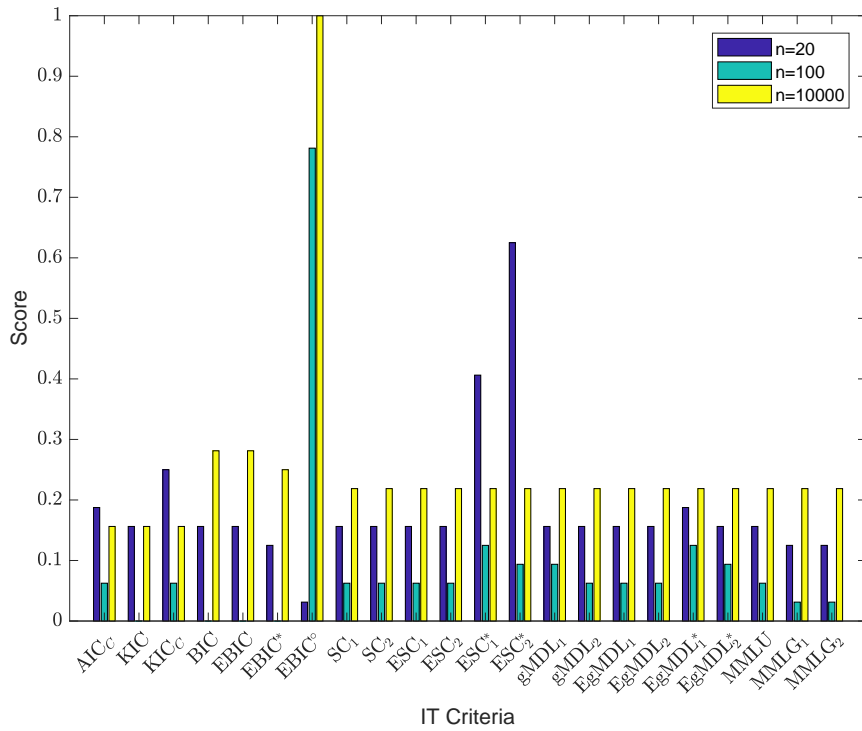


Figure 5: Scores aggregated from all experiments done for Case 3.

2 Additional information on experiments with air pollution data

2.1 Location of the four sites where the concentrations of the air pollutants are measured

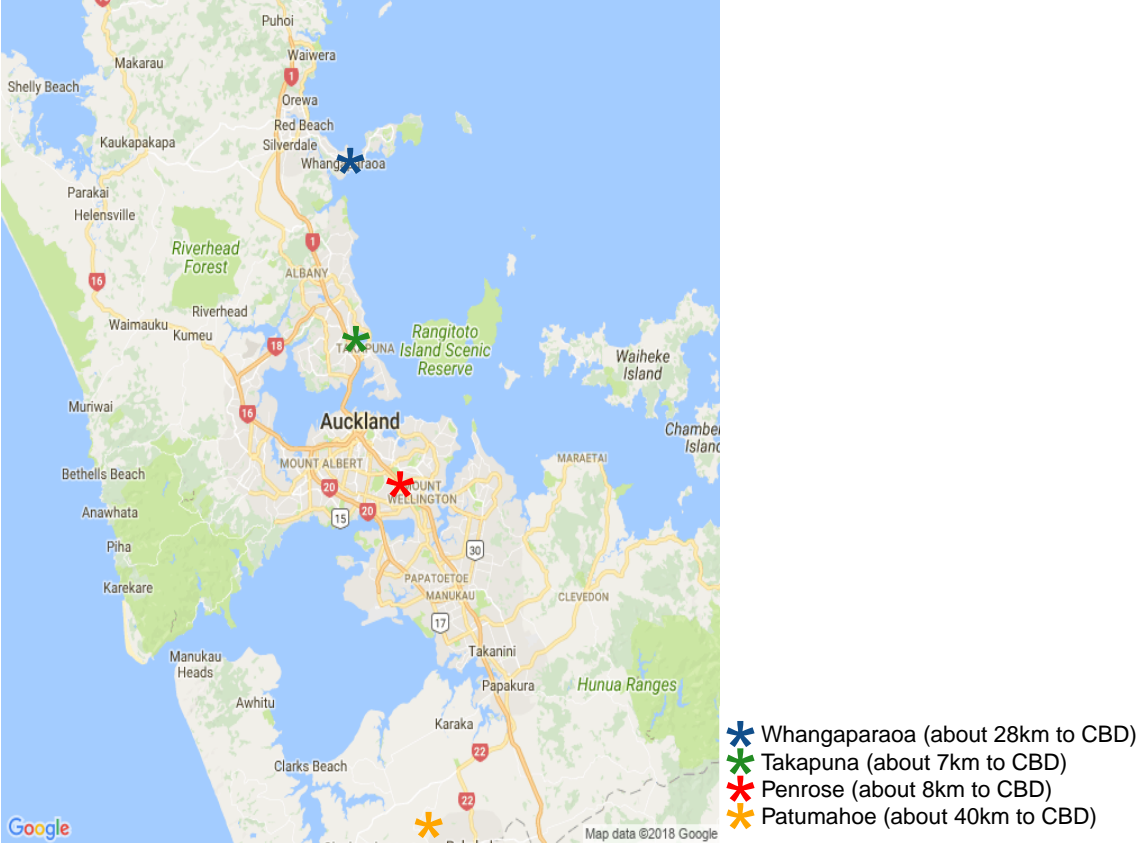


Figure 6: The distances shown in the legend are measured from Auckland central business district (CBD).

2.2 Detailed information on the experiments with air pollution data collected from Patumahoe site

2.2.1 The number of iterations and the number of predictors

In Tables 37-38 we report, for all stopping rules, the statistics computed from 100 data frames (see [1, Sec. 5.2] for the definition of the data frame). The number of iterations and the number of predictors produced by five stopping rules (AIC_C , $EBIC^\circ$, ESC_2 , ESC_2^* , CV) for each data frame are shown in Figures 7-8.

Quantiles	FullSet			ConSet		
	0.025	0.50	0.975	0.025	0.50	0.975
AIC_C	9562.0	14834.5	19998.0	96.0	239.5	487.0
KIC	20000.0	20000.0	20000.0	74.0	119.0	287.0
KIC_C	4602.0	7474.0	13474.0	71.0	113.0	246.0
BIC	19999.0	20000.0	20000.0	50.0	64.5	103.0
EBIC	19999.0	20000.0	20000.0	45.0	61.0	96.0
EBIC*	44.0	62.0	94.0	37.0	47.5	62.0
EBIC $^\circ$	27.0	32.0	40.0	32.0	37.0	44.0
SC_1	2298.0	4687.5	10849.0	48.0	61.5	128.0
SC_2	2258.0	4912.5	11677.0	48.0	61.5	128.0
ESC_1	667.0	1609.5	4148.0	44.0	59.0	104.0
ESC_2	563.0	1546.0	5588.0	44.0	59.0	104.0
ESC_1^*	43.0	62.0	94.0	37.0	47.0	62.0
ESC_2^*	41.0	59.5	94.0	37.0	47.0	62.0
$gMDL_1$	2298.0	4687.5	10849.0	49.0	66.0	128.0
$gMDL_2$	2258.0	4912.5	11677.0	49.0	66.0	128.0
$EgMDL_1$	667.0	1609.5	4148.0	45.0	60.5	117.0
$EgMDL_2$	563.0	1552.5	5588.0	45.0	60.5	110.0
$EgMDL_1^*$	44.0	63.5	94.0	37.0	47.0	62.0
$EgMDL_2^*$	44.0	62.0	94.0	37.0	47.0	62.0
MMLU	2298.0	4962.0	11677.0	48.0	61.5	126.0
$MMLG_1$	8195.0	17609.0	20000.0	45.0	61.0	136.0
$MMLG_2$	9154.0	17959.0	20000.0	45.0	61.0	136.0
CV	111.0	4875.0	20000.0	94.0	241.5	6515.0

Table 37: Summary table for the number of iterations.

Quantiles	FullSet			ConSet		
	0.025	0.50	0.975	0.025	0.50	0.975
AIC _C	513.0	590.5	654.0	18.0	27.0	36.0
KIC	623.0	646.0	678.0	14.0	20.0	29.0
KIC _C	397.0	475.5	571.0	14.0	19.0	27.0
BIC	623.0	646.0	678.0	7.0	13.0	17.0
EBIC	623.0	646.0	678.0	6.0	11.0	16.0
EBIC*	12.0	20.0	31.0	5.0	9.0	12.0
EBIC ^o	4.0	6.0	10.0	4.0	6.0	8.0
SC ₁	305.0	392.0	542.0	7.0	12.0	19.0
SC ₂	302.0	401.5	555.0	7.0	12.0	19.0
ESC ₁	149.0	252.0	375.0	6.0	11.0	17.0
ESC ₂	136.0	240.5	423.0	6.0	11.0	17.0
ESC ₁ *	12.0	20.0	32.0	5.0	8.0	12.0
ESC ₂ *	11.0	19.5	31.0	5.0	8.0	12.0
gMDL ₁	305.0	392.0	542.0	7.0	13.0	20.0
gMDL ₂	302.0	401.5	555.0	7.0	13.0	19.0
EgMDL ₁	149.0	252.0	375.0	6.0	11.0	18.0
EgMDL ₂	136.0	242.5	423.0	6.0	11.0	18.0
EgMDL ₁ *	12.0	21.0	32.0	5.0	8.0	12.0
EgMDL ₂ *	12.0	20.0	31.0	5.0	8.0	12.0
MMLU	305.0	408.0	555.0	7.0	12.0	18.0
MMLG ₁	494.0	615.5	654.0	7.0	12.0	20.0
MMLG ₂	513.0	618.0	657.0	7.0	12.0	20.0
CV	42.0	397.5	660.0	16.0	31.5	66.0

Table 38: Summary table for the number of predictors.

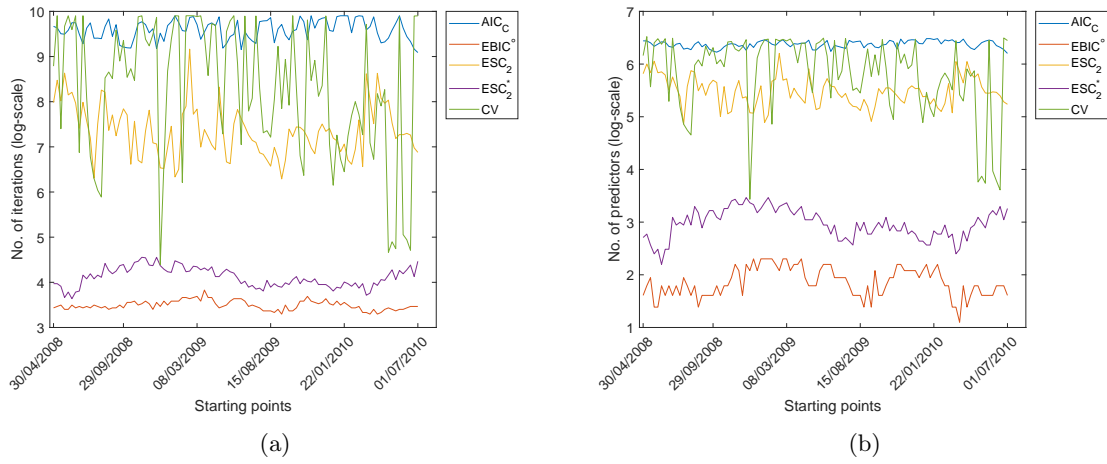


Figure 7: Results obtained for each data frame when the total number of predictors is $p_n = 1464$ (FullSet). For each plot, the starting points of the data frames are marked on the horizontal axis.

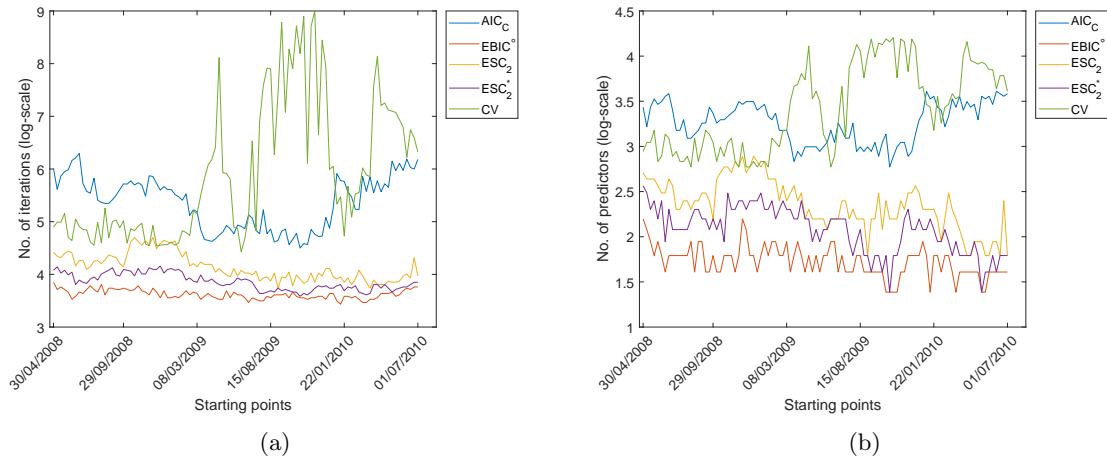


Figure 8: Results similar to those reported in Figure 7: The total number of predictors is $p_n = 68$ (ConSet).

2.2.2 Statistics on how many times each predictor is selected in 100 runs

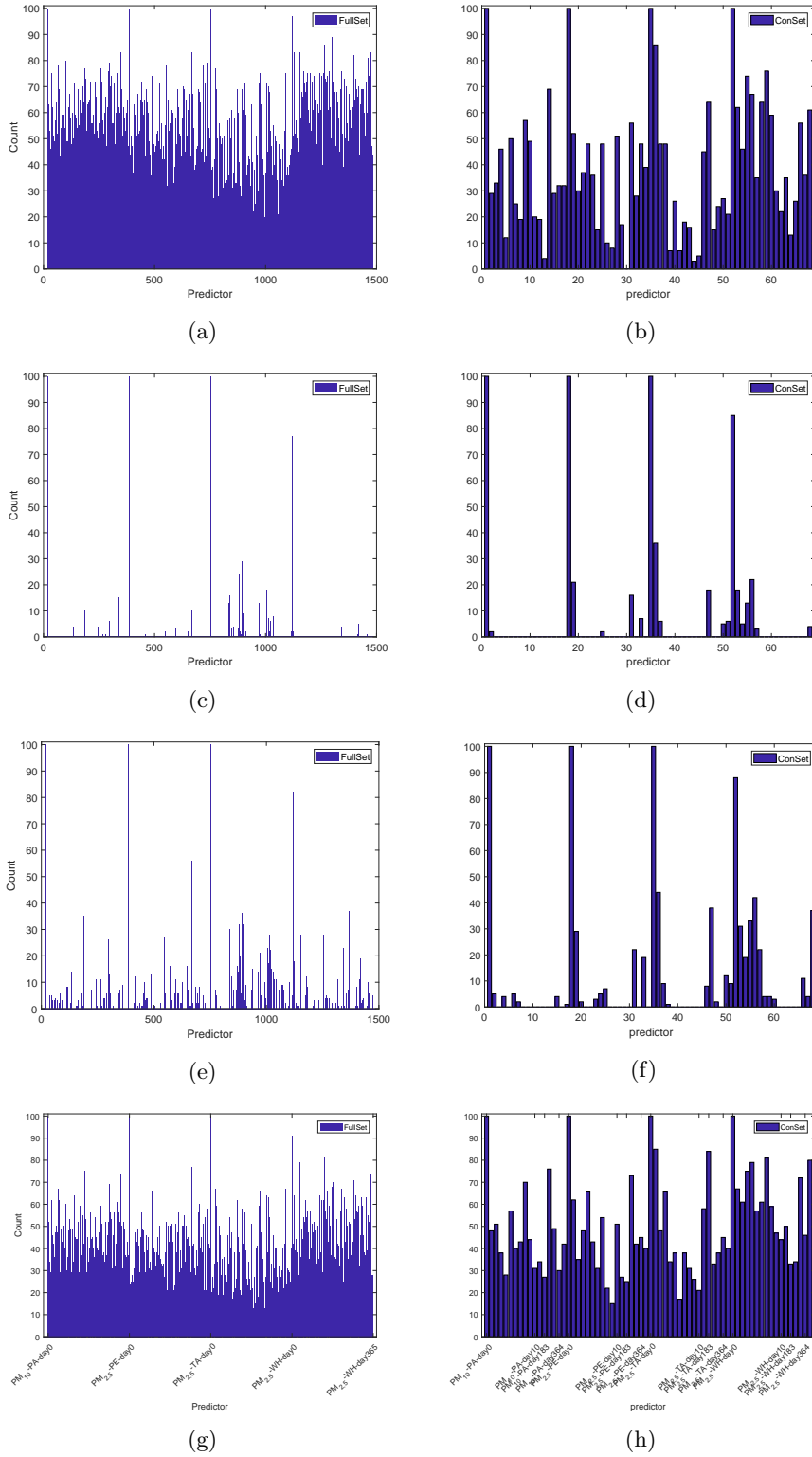


Figure 9: Results obtained with AIC_C [panels (a)-(b)], $EBIC^o$ [Panels (c)-(d)], $EgMDL_1^*$ [panels (e)-(f)] and CV [panels (g)-(h)]. On the horizontal axes, the predictors are enumerated from left to right, starting from the day t when the prediction is done and going back to day $t - 365$. Day t is marked on the graphs as day 0.

2.2.3 How many times each predictor from the constrained set is selected in 100 runs: The case when the total number of predictors is $p_n = 68$ (ConSet) versus the case when the total number of predictors is $p_n = 1464$ (FullSet)

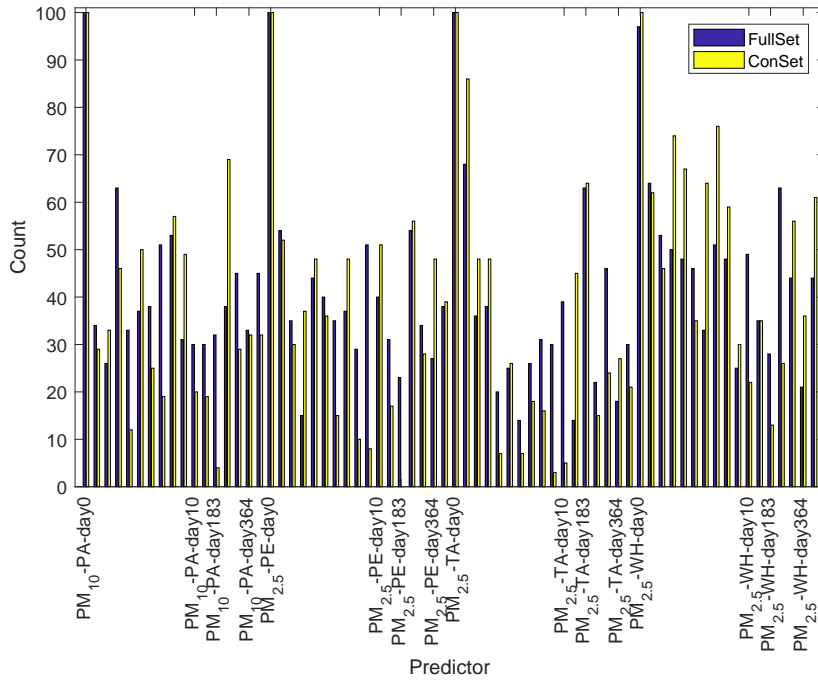


Figure 10: Results obtained when AIC_C is applied.

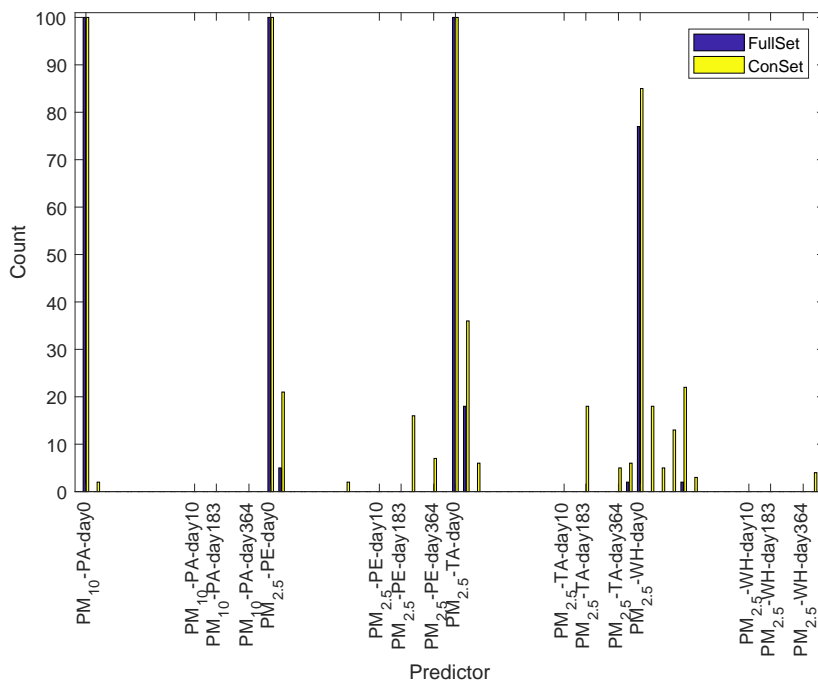


Figure 11: Results obtained when $EBIC^o$ is applied.

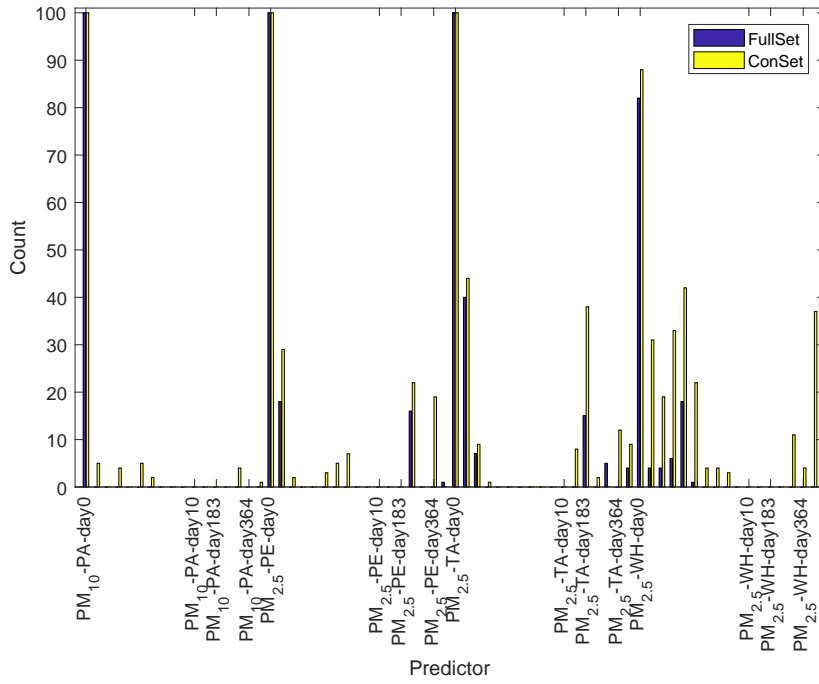


Figure 12: Results obtained when EgMDL₁* is applied.

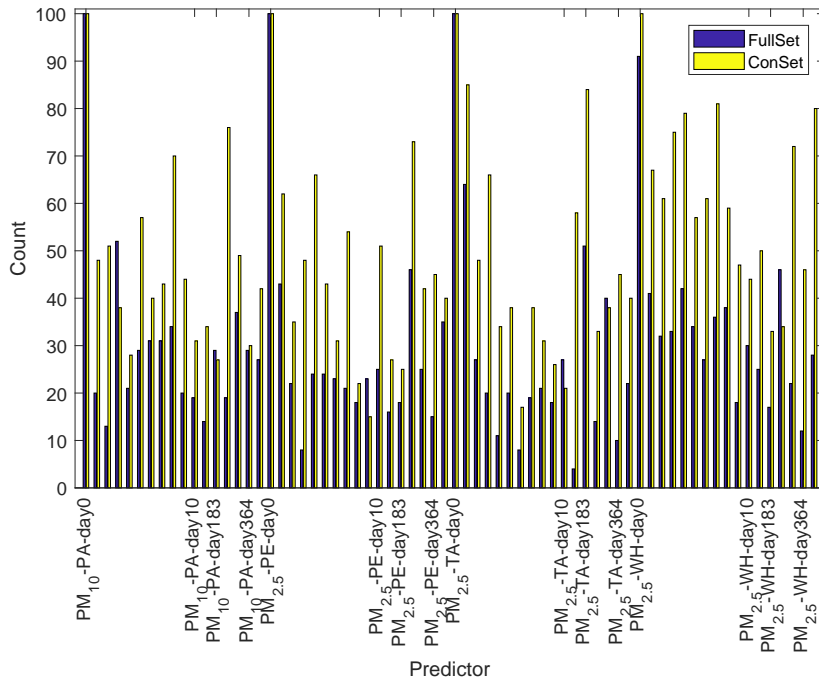


Figure 13: Results obtained when CV is applied.

2.2.4 Statistics for EgMDL_1^* on how many times each predictor from the FullSet ($p_n = 1464$) not included in the ConSet ($p_n = 68$) is selected in 100 runs

Predictor	Count
PM2.5-PE-day280	56
PM2.5-WH-day249	37
PM2.5-TA-day140	36
PM10-PA-day167	35
PM2.5-TA-day128	32
PM2.5-TA-day145	32
PM2.5-TA-day84	30
PM10-PA-day317	28
PM2.5-TA-day261	28
PM2.5-WH-day36	28
PM2.5-WH-day134	28
PM2.5-PE-day161	27
PM10-PA-day277	26
PM2.5-TA-day251	23
PM2.5-WH-day224	23
PM2.5-TA-day267	22
PM2.5-TA-day217	21
PM10-PA-day233	20
PM2.5-TA-day129	20

Table 39: List of predictors which have been selected at least twenty times in one hundred runs.

References

- [1] F. Li, C.M. Triggs, B. Dumitrescu and C.D. Giurcăneanu. The matching pursuit algorithm revisited: A variant for big data and new stopping rules. 2018.