

Worked Example for All Possible Regressions and Data Mining

This example uses the ASA software integrated into Excel or SPSS (www.asastat.com). ASA is, in part, a point-and-click interface to R but analyses can be conducted from within SPSS or Excel. All data are hypothetical. We assume you have read the primer on all possible regressions.

We use a somewhat unorthodox example because we want to illustrate the case of working with a large number of predictors where we are subject to the curse of dimensionality and to illustrate variable selection approach of Cai, Tsay and Chen (2009) described in Chapter 11. We generated hypothetical data for 63 potential predictors of an outcome variable, Y, in which the predictors were intercorrelated with one another at a magnitude of about 0.20. In the spirit of Cai et al., we initially randomly divided the predictors into 9 classes of 7 predictors each. The first set of predictors is labeled A1 through A7, the second set B1 through B7, the third set C1 through C7, and so on through I1 to I7. The sample size was 250, which is small for 63 predictors (only 4 cases per predictor). There are 2,016 different correlations in the lower triangle of the correlation matrix between all variables! We approached the classification of predictors into the sets in a completely atheoretical fashion, i.e., as if set assignment was random. In practice, theory might help assist in the classification process, but we operate throughout the exercise as if we are theoretically blind. All predictors were simulated to be normally distributed with a mean of zero and a standard deviation of 1.0.

Of the 63 predictors, only seven of them were defined so as to influence Y in accord with the following population regression equation $Y = 0.40*A1 + 0.40*A2 + 0.40*C3 + 0.40*C4 + 0.40*D5 + 0.40*D6 + 0.40*E7$. Thus, in reality, only variables A1, A2, C3, C4, D5, D6 and E7 impact Y. The overall squared R in the population for these seven predictors was set to 0.50 through the addition of normally distributed error. The influence of each of the seven predictors is moderate in magnitude, additive, and linear. Of interest is whether these seven variables are in the final set of predictors identified by the Cai et al. method as being relevant and how many false positive predictors enter the mix as well.

THE ANALYSIS

We used the ASA program called “Multiple regression: All possible regressions” in the

folder Multiple Regression – General > Multiple Regression Analyses. As a first step, we specified Y as the outcome and the first set of variables A1 through A7 as predictors. The program begins by renaming the predictors and then specifying the regression equations that were evaluated using those new names, numbering the equations from m1 to m127 (the intercept is represented by the number 1, but it is formally estimated in each model):

PREDICTOR VARIABLES

```
x1: A1  
x2: A2  
x3: A3  
x4: A4  
x5: A5  
x6: A6  
x7: A7
```

MODELS EVALUATED (1 = Intercept)

```
m1: y ~ 1 + x1 + x2 + x3 + x4 + x5 + x6 + x7  
m2: y ~ 1 + x2 + x3 + x4 + x5 + x6 + x7  
m3: y ~ 1 + x1 + x3 + x4 + x5 + x6 + x7  
m4: y ~ 1 + x3 + x4 + x5 + x6 + x7  
m5: y ~ 1 + x1 + x2 + x4 + x5 + x6 + x7  
m6: y ~ 1 + x2 + x4 + x5 + x6 + x7  
m7: y ~ 1 + x1 + x4 + x5 + x6 + x7  
m8: y ~ 1 + x4 + x5 + x6 + x7  
m9: y ~ 1 + x1 + x2 + x3 + x5 + x6 + x7  
m10: y ~ 1 + x2 + x3 + x5 + x6 + x7  
m11: y ~ 1 + x1 + x3 + x5 + x6 + x7  
m12: y ~ 1 + x3 + x5 + x6 + x7  
m13: y ~ 1 + x1 + x2 + x5 + x6 + x7  
m14: y ~ 1 + x2 + x5 + x6 + x7  
m15: y ~ 1 + x1 + x5 + x6 + x7  
m16: y ~ 1 + x5 + x6 + x7  
m17: y ~ 1 + x1 + x2 + x3 + x4 + x6 + x7  
m18: y ~ 1 + x2 + x3 + x4 + x6 + x7  
m19: y ~ 1 + x1 + x3 + x4 + x6 + x7  
m20: y ~ 1 + x3 + x4 + x6 + x7  
m21: y ~ 1 + x1 + x2 + x4 + x6 + x7  
m22: y ~ 1 + x2 + x4 + x6 + x7  
m23: y ~ 1 + x1 + x4 + x6 + x7  
m24: y ~ 1 + x4 + x6 + x7  
m25: y ~ 1 + x1 + x2 + x3 + x6 + x7  
m26: y ~ 1 + x2 + x3 + x6 + x7  
m27: y ~ 1 + x1 + x3 + x6 + x7  
m28: y ~ 1 + x3 + x6 + x7  
m29: y ~ 1 + x1 + x2 + x6 + x7
```

```
m30: y ~ 1 + x2 + x6 + x7
m31: y ~ 1 + x1 + x6 + x7
m32: y ~ 1 + x6 + x7
m33: y ~ 1 + x1 + x2 + x3 + x4 + x5 + x7
m34: y ~ 1 + x2 + x3 + x4 + x5 + x7
m35: y ~ 1 + x1 + x3 + x4 + x5 + x7
m36: y ~ 1 + x3 + x4 + x5 + x7
m37: y ~ 1 + x1 + x2 + x4 + x5 + x7
m38: y ~ 1 + x2 + x4 + x5 + x7
m39: y ~ 1 + x1 + x4 + x5 + x7
m40: y ~ 1 + x4 + x5 + x7
m41: y ~ 1 + x1 + x2 + x3 + x5 + x7
m42: y ~ 1 + x2 + x3 + x5 + x7
m43: y ~ 1 + x1 + x3 + x5 + x7
m44: y ~ 1 + x3 + x5 + x7
m45: y ~ 1 + x1 + x2 + x5 + x7
m46: y ~ 1 + x2 + x5 + x7
m47: y ~ 1 + x1 + x5 + x7
m48: y ~ 1 + x5 + x7
m49: y ~ 1 + x1 + x2 + x3 + x4 + x7
m50: y ~ 1 + x2 + x3 + x4 + x7
m51: y ~ 1 + x1 + x3 + x4 + x7
m52: y ~ 1 + x3 + x4 + x7
m53: y ~ 1 + x1 + x2 + x4 + x7
m54: y ~ 1 + x2 + x4 + x7
m55: y ~ 1 + x1 + x4 + x7
m56: y ~ 1 + x4 + x7
m57: y ~ 1 + x1 + x2 + x3 + x7
m58: y ~ 1 + x2 + x3 + x7
m59: y ~ 1 + x1 + x3 + x7
m60: y ~ 1 + x3 + x7
m61: y ~ 1 + x1 + x2 + x7
m62: y ~ 1 + x2 + x7
m63: y ~ 1 + x1 + x7
m64: y ~ 1 + x7
m65: y ~ 1 + x1 + x2 + x3 + x4 + x5 + x6
m66: y ~ 1 + x2 + x3 + x4 + x5 + x6
m67: y ~ 1 + x1 + x3 + x4 + x5 + x6
m68: y ~ 1 + x3 + x4 + x5 + x6
m69: y ~ 1 + x1 + x2 + x4 + x5 + x6
m70: y ~ 1 + x2 + x4 + x5 + x6
m71: y ~ 1 + x1 + x4 + x5 + x6
m72: y ~ 1 + x4 + x5 + x6
m73: y ~ 1 + x1 + x2 + x3 + x5 + x6
m74: y ~ 1 + x2 + x3 + x5 + x6
m75: y ~ 1 + x1 + x3 + x5 + x6
m76: y ~ 1 + x3 + x5 + x6
m77: y ~ 1 + x1 + x2 + x5 + x6
m78: y ~ 1 + x2 + x5 + x6
```

m79: $y \sim 1 + x_1 + x_5 + x_6$
m80: $y \sim 1 + x_5 + x_6$
m81: $y \sim 1 + x_1 + x_2 + x_3 + x_4 + x_6$
m82: $y \sim 1 + x_2 + x_3 + x_4 + x_6$
m83: $y \sim 1 + x_1 + x_3 + x_4 + x_6$
m84: $y \sim 1 + x_3 + x_4 + x_6$
m85: $y \sim 1 + x_1 + x_2 + x_4 + x_6$
m86: $y \sim 1 + x_2 + x_4 + x_6$
m87: $y \sim 1 + x_1 + x_4 + x_6$
m88: $y \sim 1 + x_4 + x_6$
m89: $y \sim 1 + x_1 + x_2 + x_3 + x_6$
m90: $y \sim 1 + x_2 + x_3 + x_6$
m91: $y \sim 1 + x_1 + x_3 + x_6$
m92: $y \sim 1 + x_3 + x_6$
m93: $y \sim 1 + x_1 + x_2 + x_6$
m94: $y \sim 1 + x_2 + x_6$
m95: $y \sim 1 + x_1 + x_6$
m96: $y \sim 1 + x_6$
m97: $y \sim 1 + x_1 + x_2 + x_3 + x_4 + x_5$
m98: $y \sim 1 + x_2 + x_3 + x_4 + x_5$
m99: $y \sim 1 + x_1 + x_3 + x_4 + x_5$
m100: $y \sim 1 + x_3 + x_4 + x_5$
m101: $y \sim 1 + x_1 + x_2 + x_4 + x_5$
m102: $y \sim 1 + x_2 + x_4 + x_5$
m103: $y \sim 1 + x_1 + x_4 + x_5$
m104: $y \sim 1 + x_4 + x_5$
m105: $y \sim 1 + x_1 + x_2 + x_3 + x_5$
m106: $y \sim 1 + x_2 + x_3 + x_5$
m107: $y \sim 1 + x_1 + x_3 + x_5$
m108: $y \sim 1 + x_3 + x_5$
m109: $y \sim 1 + x_1 + x_2 + x_5$
m110: $y \sim 1 + x_2 + x_5$
m111: $y \sim 1 + x_1 + x_5$
m112: $y \sim 1 + x_5$
m113: $y \sim 1 + x_1 + x_2 + x_3 + x_4$
m114: $y \sim 1 + x_2 + x_3 + x_4$
m115: $y \sim 1 + x_1 + x_3 + x_4$
m116: $y \sim 1 + x_3 + x_4$
m117: $y \sim 1 + x_1 + x_2 + x_4$
m118: $y \sim 1 + x_2 + x_4$
m119: $y \sim 1 + x_1 + x_4$
m120: $y \sim 1 + x_4$
m121: $y \sim 1 + x_1 + x_2 + x_3$
m122: $y \sim 1 + x_2 + x_3$
m123: $y \sim 1 + x_1 + x_3$
m124: $y \sim 1 + x_3$
m125: $y \sim 1 + x_1 + x_2$
m126: $y \sim 1 + x_2$
m127: $y \sim 1 + x_1$

Next, the results of each model are reported:

RESULTS

MODEL	RSQR	AdjRSQR	AIC	BIC	RMSE
1	0.331	0.311	1003.6	1035.3	1.77
2	0.270	0.252	1023.3	1051.4	1.84
3	0.270	0.252	1023.3	1051.4	1.84
4	0.205	0.189	1042.5	1067.2	1.92
5	0.324	0.307	1004.1	1032.2	1.77
6	0.250	0.235	1027.9	1052.6	1.86
7	0.257	0.242	1025.7	1050.3	1.85
8	0.174	0.160	1050.3	1071.4	1.95
9	0.328	0.311	1002.7	1030.9	1.77
10	0.263	0.248	1023.6	1048.3	1.84
11	0.262	0.247	1024.1	1048.8	1.85
12	0.190	0.177	1045.3	1066.4	1.93
13	0.320	0.306	1003.5	1028.1	1.77
14	0.241	0.228	1029.2	1050.3	1.87
15	0.246	0.234	1027.3	1048.4	1.86
16	0.152	0.142	1054.8	1072.4	1.97
17	0.329	0.312	1002.3	1030.5	1.76
18	0.266	0.251	1022.8	1047.5	1.84
19	0.262	0.247	1023.9	1048.6	1.85
20	0.192	0.179	1044.7	1065.8	1.93
21	0.322	0.308	1003.0	1027.6	1.77
22	0.244	0.231	1028.2	1049.3	1.87
23	0.247	0.235	1027.1	1048.3	1.86
24	0.154	0.144	1054.2	1071.8	1.97
25	0.326	0.312	1001.5	1026.2	1.77
26	0.258	0.246	1023.5	1044.7	1.85
27	0.252	0.240	1025.4	1046.5	1.86
28	0.173	0.163	1048.5	1066.1	1.95
29	0.317	0.306	1002.6	1023.7	1.77
30	0.232	0.222	1030.1	1047.7	1.88
31	0.233	0.224	1029.6	1047.2	1.87
32	0.125	0.118	1060.6	1074.6	2.00
33	0.324	0.307	1004.1	1032.3	1.77
34	0.258	0.243	1025.5	1050.2	1.85
35	0.256	0.241	1026.1	1050.7	1.85
36	0.182	0.169	1047.6	1068.8	1.94
37	0.317	0.303	1004.8	1029.4	1.78
38	0.236	0.223	1030.9	1052.0	1.88
39	0.241	0.229	1029.0	1050.2	1.87
40	0.145	0.135	1056.8	1074.4	1.98
41	0.319	0.306	1003.8	1028.4	1.77
42	0.247	0.235	1027.0	1048.1	1.86
43	0.243	0.231	1028.4	1049.5	1.87

44	0.159	0.149	1052.8	1070.4	1.96
45	0.311	0.300	1004.9	1026.0	1.78
46	0.221	0.211	1033.7	1051.3	1.89
47	0.225	0.215	1032.4	1050.0	1.88
48	0.111	0.103	1064.7	1078.8	2.01
49	0.321	0.307	1003.2	1027.8	1.77
50	0.251	0.239	1025.8	1047.0	1.86
51	0.244	0.232	1028.0	1049.1	1.86
52	0.162	0.152	1051.8	1069.4	1.96
53	0.313	0.302	1004.2	1025.3	1.78
54	0.225	0.216	1032.2	1049.8	1.88
55	0.226	0.216	1032.0	1049.7	1.88
56	0.115	0.107	1063.6	1077.7	2.01
57	0.316	0.305	1003.1	1024.3	1.77
58	0.238	0.229	1028.0	1045.6	1.87
59	0.228	0.218	1031.4	1049.0	1.88
60	0.130	0.123	1059.1	1073.2	1.99
61	0.306	0.298	1004.7	1022.3	1.78
62	0.207	0.200	1036.2	1050.2	1.90
63	0.204	0.197	1037.1	1051.2	1.91
64	0.065	0.061	1075.1	1085.7	2.06
65	0.319	0.302	1006.1	1034.3	1.78
66	0.259	0.244	1025.0	1049.7	1.85
67	0.250	0.235	1028.1	1052.8	1.86
68	0.186	0.173	1046.5	1067.6	1.93
69	0.313	0.299	1006.3	1030.9	1.78
70	0.241	0.228	1029.1	1050.3	1.87
71	0.238	0.225	1030.1	1051.3	1.87
72	0.156	0.146	1053.5	1071.1	1.97
73	0.316	0.302	1005.1	1029.7	1.78
74	0.253	0.241	1025.2	1046.3	1.85
75	0.242	0.230	1028.8	1049.9	1.87
76	0.172	0.162	1048.9	1066.5	1.95
77	0.309	0.298	1005.5	1026.6	1.78
78	0.232	0.222	1030.1	1047.7	1.88
79	0.228	0.218	1031.5	1049.1	1.88
80	0.135	0.128	1057.7	1071.8	1.99
81	0.314	0.300	1005.6	1030.3	1.78
82	0.251	0.239	1025.6	1046.8	1.86
83	0.235	0.223	1031.0	1052.1	1.88
84	0.164	0.154	1051.1	1068.7	1.96
85	0.308	0.296	1006.1	1027.2	1.78
86	0.230	0.221	1030.6	1048.2	1.88
87	0.220	0.210	1034.0	1051.6	1.89
88	0.126	0.119	1060.2	1074.3	2.00
89	0.311	0.300	1004.8	1025.9	1.78
90	0.244	0.235	1026.2	1043.8	1.86
91	0.225	0.215	1032.4	1050.0	1.89
92	0.145	0.138	1054.8	1068.9	1.98

93	0.304	0.295	1005.6	1023.2	1.79
94	0.219	0.213	1032.2	1046.3	1.89
95	0.206	0.199	1036.4	1050.5	1.90
96	0.097	0.094	1066.4	1077.0	2.03
97	0.306	0.292	1008.5	1033.2	1.79
98	0.240	0.227	1029.5	1050.6	1.87
99	0.224	0.212	1034.5	1055.6	1.89
100	0.149	0.139	1055.5	1073.1	1.97
101	0.300	0.288	1008.9	1030.0	1.79
102	0.219	0.209	1034.3	1051.9	1.89
103	0.210	0.200	1037.1	1054.7	1.90
104	0.113	0.106	1064.1	1078.2	2.01
105	0.302	0.290	1008.2	1029.3	1.79
106	0.229	0.220	1030.9	1048.5	1.88
107	0.210	0.201	1037.0	1054.6	1.90
108	0.124	0.117	1060.9	1075.0	2.00
109	0.294	0.286	1008.9	1026.5	1.80
110	0.204	0.198	1036.9	1051.0	1.91
111	0.192	0.186	1040.7	1054.7	1.92
112	0.076	0.072	1072.2	1082.8	2.05
113	0.299	0.288	1009.2	1030.3	1.80
114	0.226	0.217	1031.9	1049.5	1.88
115	0.199	0.189	1040.7	1058.3	1.92
116	0.110	0.103	1064.8	1078.9	2.02
117	0.291	0.283	1010.0	1027.6	1.80
118	0.201	0.194	1037.9	1052.0	1.91
119	0.179	0.172	1044.8	1058.9	1.94
120	0.059	0.055	1076.8	1087.4	2.07
121	0.293	0.285	1009.2	1026.8	1.80
122	0.213	0.206	1034.2	1048.3	1.90
123	0.178	0.171	1045.0	1059.1	1.94
124	0.072	0.068	1073.4	1084.0	2.05
125	0.284	0.278	1010.5	1024.5	1.81
126	0.181	0.178	1042.0	1052.6	1.93
127	0.152	0.148	1050.9	1061.5	1.96

We illustrate the approach to model selection that uses the BIC. First, we locate the model with the lowest BIC, which is model 61 with a BIC of 1022.3. This is the initial “best” model. The predictors in this model are A1, A2 and A7. Raftery (1995) suggests that if two models are within 2.2 BIC units of each other, they have about equal support. We therefore scan the list to determine if there is a more parsimonious model that is within 2.2 BIC units of 1022.3. Model 125 has 2 predictors, A1 and A2, and its BIC is 2.2 units larger than model 61. We therefore used this model, selecting A1 and A2 from the first set as the “winning” predictors for our next iteration, described later.

We repeated this process for each of the remaining 8 sets of predictors, one set at a time. In cases where the number of predictors were the same in the most parsimonious

equations that were within 2.2 BIC units of the lowest BIC value, we selected the equation with the lowest BIC. The predictors that entered the final winning pool after all nine sets were analyzed were A1, A2, B4, B5, B7, C3, C4, C6, D3, D5, D6, E1, E2, E7, F2, G1, H2 and I7, or 18 predictors. We randomly assigned these predictors to three sets of 6 predictors each and repeated the BIC model selection process for each of these three sets, separately, just as we did in the first iteration. The winning predictors from these three analyses were A1, A2, C3, C4, D5, D6, D7, E2 and E7, or only 9 predictors. We decided to submit all 9 predictors to an all possible regression analysis, with the final selected model that had the lowest BIC having the predictors A1, A2, C3, C4, D5, D6 and E7. These are the seven predictors that were indeed the population determinants of Y! One might then consider theory building around these seven variables.

Although the approach was effective in this case, there undoubtedly will be scenarios where it includes false positives or false negatives. Again, we caution against approaching theory construction in a theory blind fashion.

REFERENCES

- Cai, A., Tsay, R. & Chen, R. (2009). Variable selection in linear regression with many predictors, *Journal of Computational and Graphical Statistics*, 18, 573-591.
- Raftery, A.E. (1995). Bayesian model selection in social research (with discussion). *Sociological Methodology*, 25, 111-195.