

**Stat 5100 Handout #27 – SAS: Variations on Ordinary Least Squares  
(LASSO and Elastic Net)**

Example: (Baseball) This data set (from the SAS Help) contains salary (for 1987) and performance (1986 and some career) data for 322 MLB players who played at least one game in both 1986 and 1987 seasons, excluding pitchers. How can salary be predicted from performance?

```
data baseball; set sashelp.baseball;
proc contents varnum data=baseball;
ods select position;
run;
```

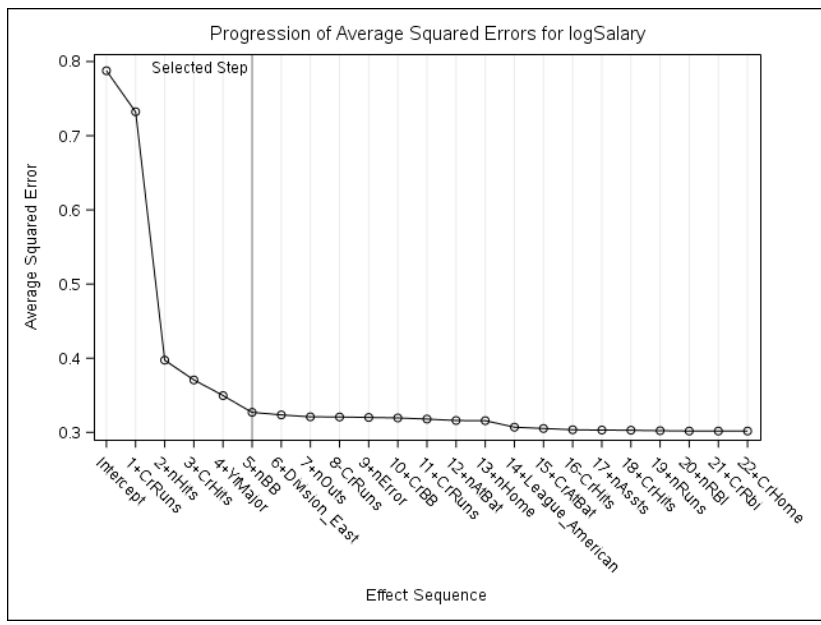
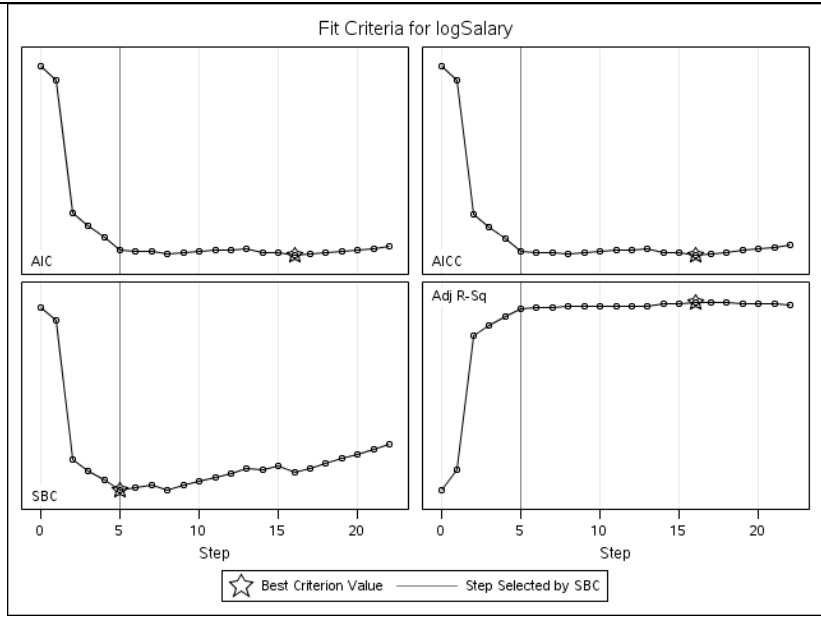
Variables in Creation Order				
#	Variable	Type	Len	Label
1	Name	Char	18	Player's Name
2	Team	Char	14	Team at the End of 1986
3	nAtBat	Num	8	Times at Bat in 1986
4	nHits	Num	8	Hits in 1986
5	nHome	Num	8	Home Runs in 1986
6	nRuns	Num	8	Runs in 1986
7	nRBI	Num	8	RBI's in 1986
8	nBB	Num	8	Walks in 1986
9	YrMajor	Num	8	Years in the Major Leagues
10	CrAtBat	Num	8	Career Times at Bat
11	CrHits	Num	8	Career Hits
12	CrHome	Num	8	Career Home Runs
13	CrRuns	Num	8	Career Runs
14	CrRbi	Num	8	Career RBI's
15	CrBB	Num	8	Career Walks
16	League	Char	8	League at the End of 1986
17	Division	Char	8	Division at the End of 1986
18	Position	Char	8	Position(s) in 1986
19	nOuts	Num	8	Put Outs in 1986
20	nAssts	Num	8	Assists in 1986
21	nError	Num	8	Errors in 1986
22	Salary	Num	8	1987 Salary in \$ Thousands
23	Div	Char	16	League and Division
24	logSalary	Num	8	Log Salary

```

/* lasso */
proc glmselect data=baseball plots=(criterion ase);
class league division;
model logSalary = nAtBat nHits nHome nRuns nRBI nBB
              yrMajor crAtBat crHits crHome crRuns crRbi
              crBB league division nOuts nAssts nError
  / selection=lasso(adaptive choose=sbc stop=none);
output out=out1 p=predlasso;
run;

```

Data Set			LASSO Selection Summary				
Dependent Variable	logSalary		Step	Effect Entered	Effect Removed	Number Effects In	SBC
Selection Method	Adaptive LASSO		* Optimal Value of Criterion				
Stop Criterion	None		0	Intercept		1	-57.2041
Choose Criterion	SBC		1	CrRuns		2	-70.8348
Effect Hierarchy Enforced	None		2	nHits		3	-226.0696
<hr/>			3	CrHits		4	-238.6648
<hr/>			4	YrMajor		5	-248.4971
<hr/>			5	nBB		6	-260.5682*
<hr/>			6	Division_East		7	-257.7020
<hr/>			7	nOuts		8	-254.3352
<hr/>			8		CrRuns	7	-260.1040
<hr/>			9	nError		8	-254.9990
<hr/>			10	CrBB		9	-249.9243
<hr/>			11	CrRuns		10	-245.7008
<hr/>			12	nAtBat		11	-241.6564
<hr/>			13	nHome		12	-236.3245
<hr/>			14	League_American		13	-238.1068
<hr/>			15	CrAtBat		14	-234.0015
<hr/>			16		CrHits	13	-241.0870
<hr/>			17	nAssts		14	-235.9894
<hr/>			18	CrHits		15	-230.5456
<hr/>			19	nRuns		16	-225.5197
<hr/>			20	nRBI		17	-220.3634
<hr/>			21	CrRbi		18	-214.7952
<hr/>			22	CrHome		19	-209.2505
<hr/>			Selection stopped because all candidate effects for entry are linearly dependent on effects in the model.				



**Selected Model**  
**The selected model, based on SBC, is the model at Step 5.**

Root MSE	0.57845
Dependent Mean	5.92722
R-Square	0.5849
Adj R-Sq	0.5768
AIC	-17.00115
AICC	-16.56194
SBC	-260.56823

Parameter Estimates		
Parameter	DF	Estimate
Intercept	1	4.229778
nHits	1	0.007194
nBB	1	0.005629
YrMajor	1	0.062808
CrHits	1	0.000222
CrRuns	1	0.000136

```

/* elastic net */
proc glmselect data=out1 plots=(criterion ase) seed=12;
class league division;
model logSalary = nAtBat nHits nHome nRuns nRBI nBB
              yrMajor crAtBat crHits crHome crRuns crRbi
              crBB league division nOuts nAssts nError
/ selection=elasticnet(stop=none choose=cv)
  cvmethod=random(20);
output out=out2 p=predelasticnet;
run;

```

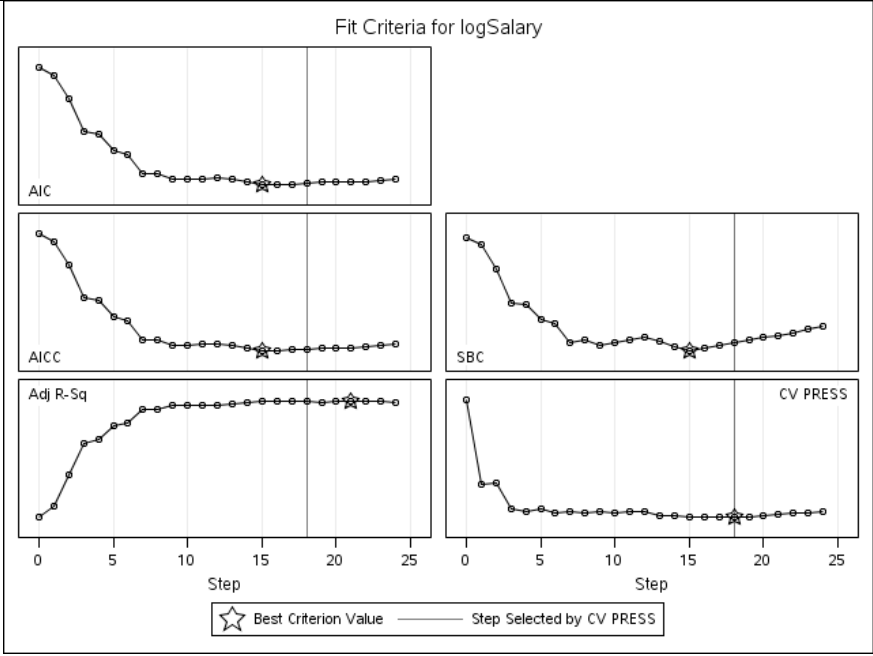
Data Set			Elastic Net Selection Summary				
Dependent Variable				Effect	Effect	Number	CV
Selection Method			Step	Entered	Removed	Effects	PRESS
WORK.OUT1							
logSalary							
ELASTICNET							
None							
Cross Validation			* Optimal Value of Criterion				
Random			0	Intercept		1	209.2326
20			1	CrRuns		2	123.1776
None			2	CrHits		3	123.7433
12			3	nHits		4	97.6956
			4	nBB		5	94.7216
			5	CrRbi		6	98.1015
			6	YrMajor		7	92.7082
			7	nRBI		8	94.5500
			8	Division_East		9	93.3921
			9	nOuts		10	94.1530
			10	nError		11	93.8913
			11	nHome		12	94.2533
			12	League_American		13	94.4968
			13		CrRbi	12	90.7314
			14		CrRuns	11	90.1957
			15		nRBI	10	89.6571
			16	CrBB		11	89.2733
			17	CrRuns		12	89.4515
			18	nAtBat		13	88.9017*
			19	CrAtBat		14	89.2818
			20	nAssts		15	89.7926
			21	CrHome		16	91.8598
			22	nRBI		17	92.6309
			23	nRuns		18	93.1973
			24	CrRbi		19	94.5881

Class Level Information		
Class	Levels	Values
League	2	American National
Division	2	East West

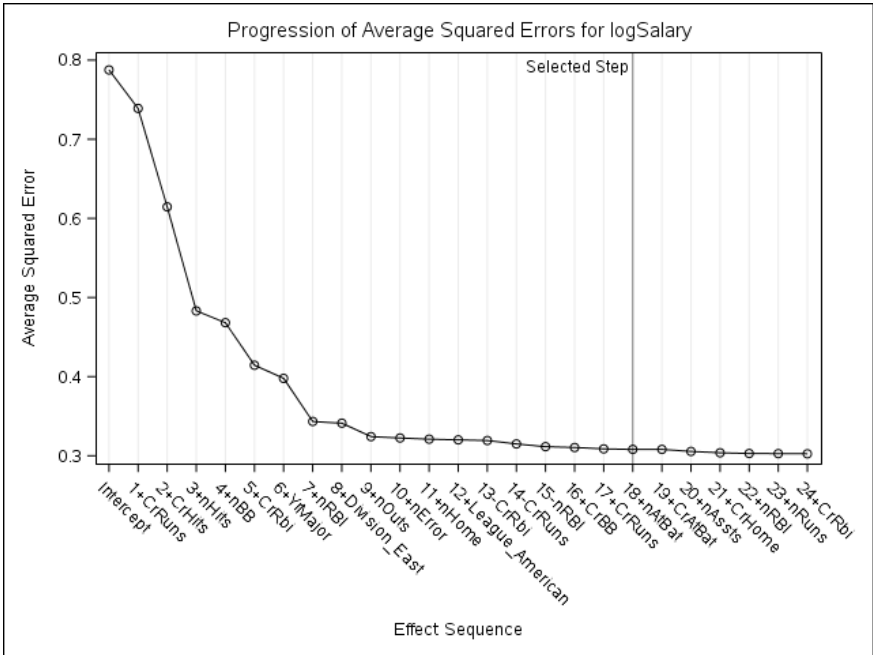
\* Optimal Value of Criterion

Selection stopped because all candidate effects for entry are linearly dependent on effects in the model.



**Selected Model**  
**The selected model, based on Cross Validation, is the model at Step 18.**

Root MSE	0.56923
Dependent Mean	5.92722
R-Square	0.6090
Adj R-Sq	0.5902
AIC	-18.72037
AICC	-17.02682
SBC	-237.28237
CV PRESS	88.90168



**Parameter Estimates**

Parameter	D	F	Estimate
Intercept	1		4.195962
nAtBat	1		-0.000112
nHits	1		0.006807
nHome	1		0.003545
nBB	1		0.007082
YrMajor	1		0.070194
CrHits	1		0.000247
CrRuns	1		0.000212
CrBB	1		-0.000348
League_American	1		-0.092575
Division_East	1		0.144062
nOurs	1		0.000192
nError	1		-0.007767

```
proc sgscatter data=out2;
  matrix logSalary predlasso predelasticnet /
    markerattrs=(symbol=circlefilled size=6pt);
run;
```

