

# STOR 455

# **STATISTICAL METHODS I**

Jan Hannig

# Model Selection (Chapter 7)

- Model selection methods
  - All subset selection
  - Forward selection, backward elimination, stepwise regression

# All-subset selection

- Consider all subset of the full model
- Select submodel based on certain criterion
- May have several “good models” instead of one “best model”
- Number of models to consider is  $2^p$ , not possible when  $p$  is large!

# Criteria for Variable Selection

- Variance based: **minimize root-MSE**
- SSE based criterion: **maximize**  
 $R^2 = 1 - \text{SSE}(\text{model}) / \text{SSTO}$
- MSE based criterion: **maximize**  
**adjusted- $R^2$**   $= 1 - \text{MSE}(\text{model}) / \text{MSTO}$   
accounting for d.f. of the model
- **Mallow's  $C_p$** : minimizing the bias of sub-model
- **PRESS<sub>p</sub>**: minimizing prediction error

# Mallow's $C_p$

- Compare subset models with the full model
  - A subset model is good if there is not substantial bias in the predicted values (relative to the full model)
- $C_p$  is a measure of this bias
$$C_p = SSE_p / MSE_{Full} - (n - 2p)$$
- If no bias  $C_p \approx p$ , therefore we prefer models with smaller  $C_p$ .

# PRESS

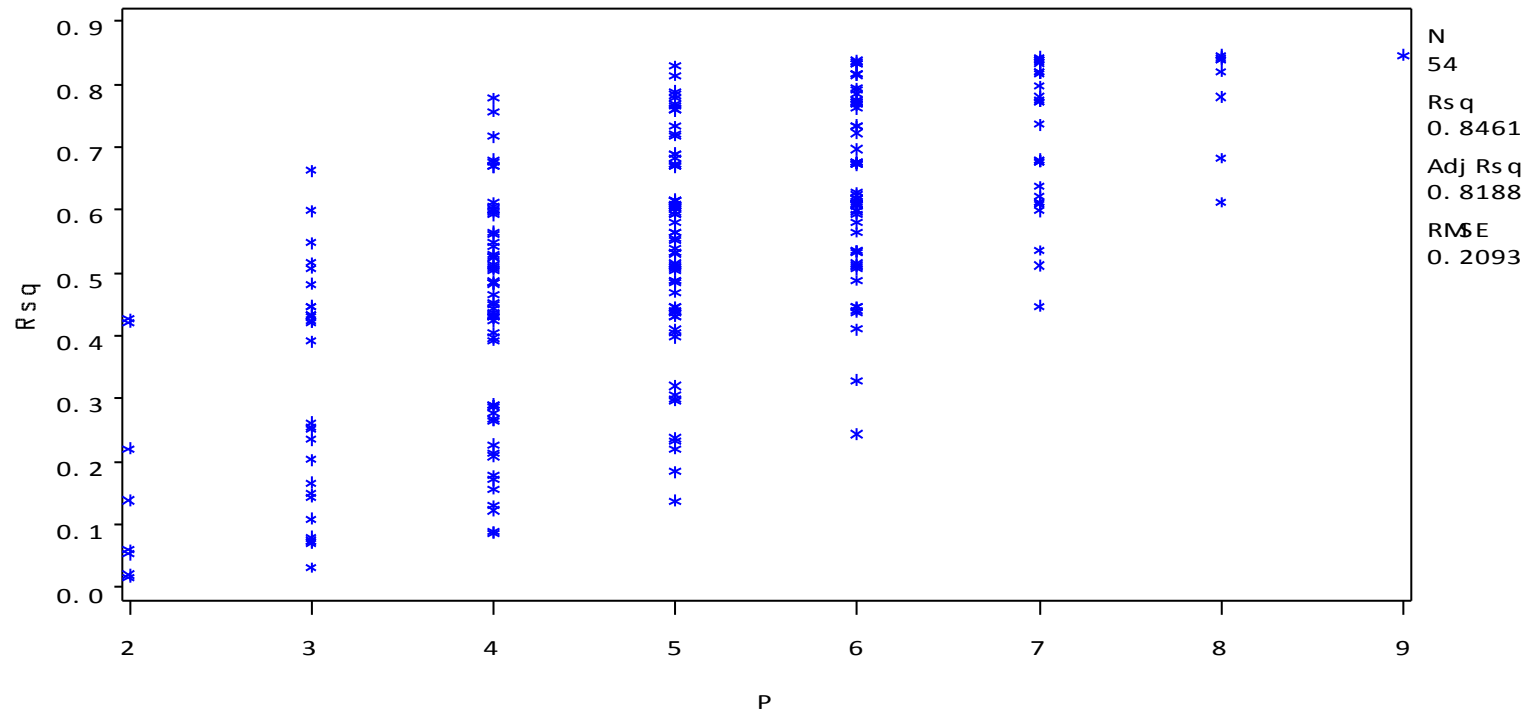
- Prediction error sums of squares
- For each case  $i$ , delete the case and predict  $Y$  using a model based on the other  $n-1$  cases
- $\text{PRESS} = \text{SS for observed minus predicted}$
- Small PRESS indicate better fit

# Do it in SAS

```
* All subset procedure, R^2, adj. R^2 and cp  
  as criteria;  
* Plot may not be informative;  
options reset=all;  
symbol1 v=star c=blue h = .8;  
proc reg data = surg;  
  model lsurv = blood prog enz liver age  
  gend alc1 alc2/selection= rsquare adjrsq  
  cp mse;  
  plot rsq.*np.;  
  plot mse.*np.;  
  plot cp.*np. ;  
run;
```

# Do it in SAS

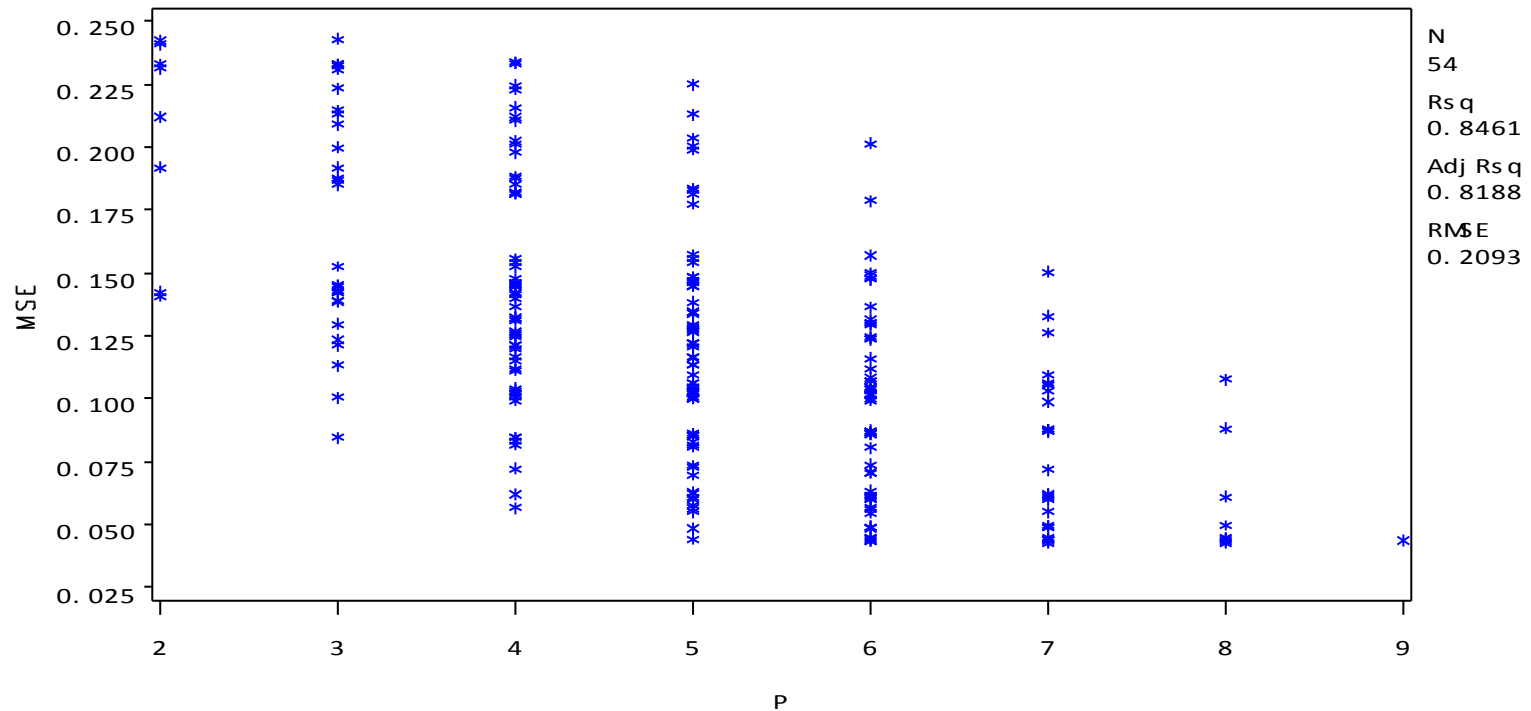
l surv = 4.0505 +0.0685 bl ood +0.0135 prog +0.015 enz +0.008 l i ver -0.0036 age +0.0842 gend  
+0.0579 al c1 +0.3884 al c2





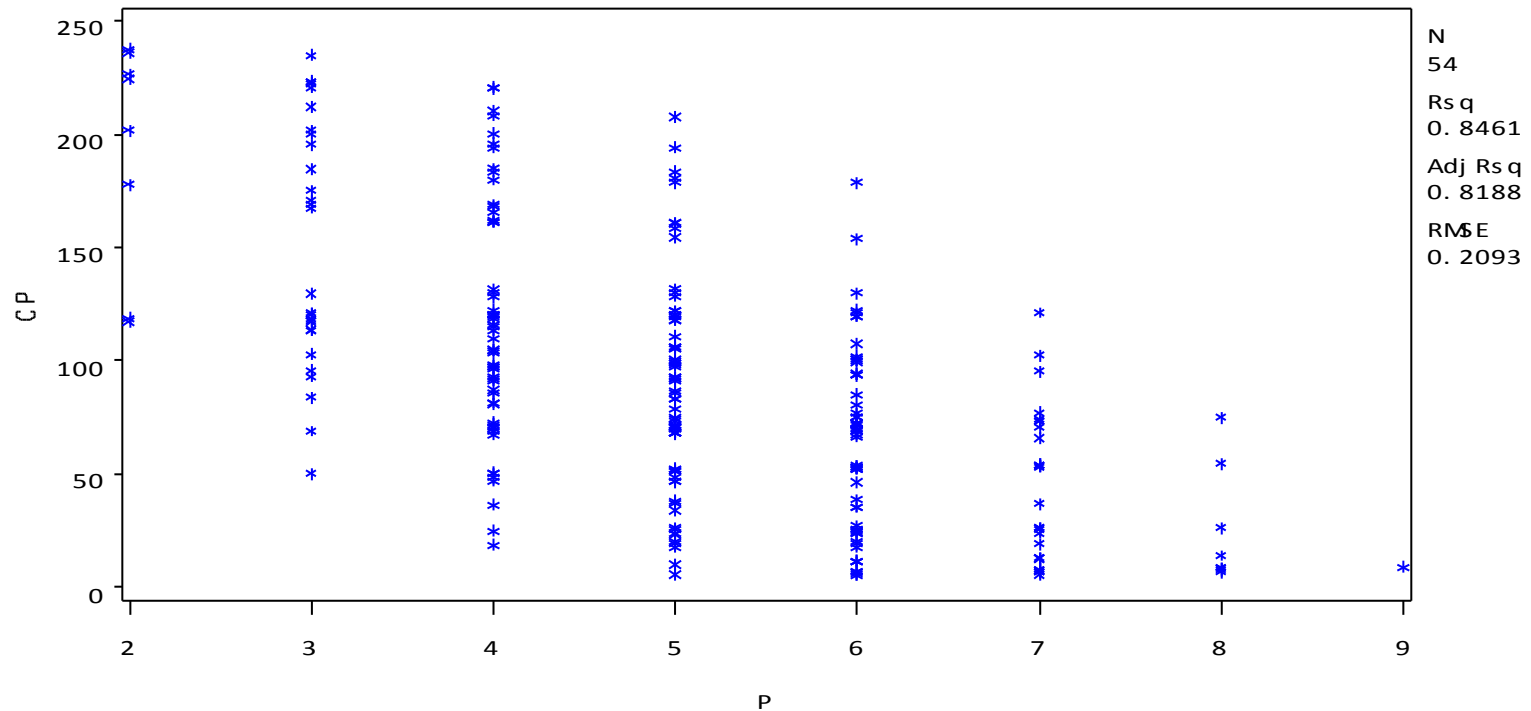
# Do it in SAS

l surv = 4.0505 +0.0685 bl ood +0.0135 prog +0.015 enz +0.008 l i ver -0.0036 age +0.0842 gend  
+0.0579 al c1 +0.3884 al c2



# Do it in SAS

l surv = 4.0505 +0.0685 bl ood +0.0135 prog +0.015 enz +0.008 l i ver -0.0036 age +0.0842 gend  
+0.0579 al c1 +0.3884 al c2



# Do it in SAS

## R-Square Selection Method

Number in Model	R-Square	Adjusted R-Square	C(p)	MSE	Variables in Model
3	0.7780	0.7647	18.9145	0.05686	prog enz alc2
3	0.7573	0.7427	24.9805	0.06217	blood prog enz
3	0.7178	0.7009	36.5252	0.07228	prog enz liver
3	0.6810	0.6618	47.3038	0.08172	prog enz alc1
-----					
4	0.8299	0.8160	5.7508	0.04447	blood prog enz alc2
4	0.8144	0.7993	10.2670	0.04850	prog enz liver alc2
4	0.7888	0.7715	17.7770	0.05521	prog enz gend alc2
4	0.7836	0.7659	19.2976	0.05657	prog enz age alc2
-----					
5	0.8374	0.8205	5.5406	0.04338	blood prog enz gend alc2
5	0.8358	0.8187	6.0182	0.04381	blood prog enz age alc2
5	0.8331	0.8158	6.7986	0.04452	blood prog enz liver alc2
5	0.8317	0.8141	7.2269	0.04491	blood prog enz alc1 alc2
5	0.8179	0.7989	11.2608	0.04859	prog enz liver gend alc2
5	0.8159	0.7968	11.8266	0.04911	prog enz liver alc1 alc2
5	0.8157	0.7965	11.9099	0.04919	prog enz liver age alc2
5	0.7944	0.7730	18.1324	0.05486	prog enz age gend alc2

# Do it in SAS

## R-Square Selection Method

Number in Model	R-Square	Adjusted R-Square	C(p)	MSE	Variables in Model
-----					
6	0.8434	0.8234	5.7874	0.04266	blood prog enz age gend alc2
6	0.8392	0.8187	7.0295	0.04382	blood prog enz gend alc1 alc2
6	0.8387	0.8181	7.1662	0.04395	blood prog enz liver gend alc2
6	0.8384	0.8178	7.2462	0.04402	blood prog enz age alc1 alc2
6	0.8371	0.8164	7.6270	0.04438	blood prog enz liver age alc2
6	0.8348	0.8137	8.3146	0.04502	blood prog enz liver alc1 alc2
-----					
7	0.8460	0.8226	7.0295	0.04287	blood prog enz age gend alc1 alc2
7	0.8436	0.8198	7.7352	0.04354	blood prog enz liver age gend alc2
7	0.8404	0.8161	8.6793	0.04444	blood prog enz liver gend alc1 alc2
7	0.8396	0.8151	8.9214	0.04467	blood prog enz liver age alc1 alc2
7	0.8213	0.7941	14.2632	0.04976	prog enz liver age gend alc1 alc2
7	0.7801	0.7466	26.3115	0.06123	blood prog enz liver age gend alc1
7	0.6829	0.6347	54.7334	0.08829	blood enz liver age gend alc1 alc2
7	0.6126	0.5537	75.2932	0.10786	blood prog liver age gend alc1 alc2
-----					
8	0.8461	0.8188	9.0000	0.04379	blood prog enz liver age gend alc1 alc2

# Best Models Using of $C_p$

- blood prog enz gend alc2, 5.54
- blood prog enz alc2, 5.75
- blood prog enz age gend alc2, 5.79
- blood prog enz age gend alc1 alc2, 7.03
- blood prog enz liver age gend alc2, 7.74

# SAS MACRO allsubsreg

```
* Use macro to compute PRESS;  
%include "T:\...\allsubsreg.sas";  
%allsubsreg(data=surg, depvar=lsurv,  
  indepvar=blood prog enz liver age  
  gend alc1 alc2, sortvar=_PRESS_,  
  printvar=_RMSE_ _RSQ_ _CP_ _PRESS_);  
run;
```

# Do it in SAS

----- Number of regressors in model=4 -----  
(continued)

Obs	VarsInModel	_IN_	_RMSE_	_RSQ_	_CP_	_PRESS_	Intercept	blood
156	blood prog enz alc1	4	0.24705	0.76650	24.2885	3.88900	3.83207	0.09176
157	blood prog enz age	4	0.24579	0.76888	23.5924	3.86305	4.02810	0.09483
158	prog enz age alc2	4	0.23785	0.78356	19.2976	3.53115	4.47171	.
159	prog enz alc1 alc2	4	0.23982	0.77996	20.3519	3.52287	4.26344	.
160	prog enz gend alc2	4	0.23498	0.78876	17.7770	3.50513	4.29334	.
161	prog enz liver alc2	4	0.22023	0.81444	10.2670	3.02103	4.34067	.
162	blood prog enz alc2	4	0.21087	0.82988	5.7508	2.73777	3.85242	0.07332

----- Number of regressors in model=5 -----  
(continued)

Obs	VarsInModel	_IN_	_RMSE_	_RSQ_	_CP_	_PRESS_	Intercept
209	prog enz age alc1 alc2	5	0.23877	0.78633	20.4874	3.56560	4.45344
210	prog enz age gend alc2	5	0.23423	0.79439	18.1324	3.55562	4.47592
211	prog enz gend alc1 alc2	5	0.23636	0.79063	19.2313	3.55539	4.26651
212	prog enz liver age alc2	5	0.22178	0.81566	11.9099	3.13206	4.42572
213	prog enz liver gend alc2	5	0.22044	0.81788	11.2608	3.12603	4.33834
214	prog enz liver alc1 alc2	5	0.22161	0.81595	11.8266	3.10027	4.31629
215	blood prog enz liver alc2	5	0.21100	0.83314	6.7986	2.82935	3.96517
216	blood prog enz alc1 alc2	5	0.21193	0.83168	7.2269	2.79565	3.82669
217	blood prog enz gend alc2	5	0.20827	0.83744	5.5406	2.78271	3.86710
218	blood prog enz age alc2	5	0.20931	0.83581	6.0182	2.73893	4.03812

# Do it in SAS

----- Number of regressors in model=6 -----  
(continued)

Obs	VarsInModel	_IN_	_RMSE_	_RSQ_	_CP_	_PRESS_	Intercept
239	prog enz liver gend alc1 alc2	6	0.22185	0.81939	12.8203	3.20786	4.31396
240	prog enz liver age alc1 alc2	6	0.22298	0.81754	13.3595	3.20039	4.41122
241	blood prog enz liver alc1 alc2	6	0.21218	0.83479	8.3146	2.89906	3.93809
242	blood prog enz liver gend alc2	6	0.20964	0.83872	7.1662	2.87494	3.93868
243	blood prog enz liver age alc2	6	0.21066	0.83715	7.6270	2.85668	4.08721
244	blood prog enz gend alc1 alc2	6	0.20934	0.83919	7.0295	2.83917	3.84163
245	blood prog enz age alc1 alc2	6	0.20982	0.83845	7.2462	2.77826	4.02082
246	blood prog enz age gend alc2	6	0.20655	0.84344	5.7874	2.77233	4.05397

----- Number of regressors in model=7 -----

Obs	VarsInModel	_IN_	_RMSE_	_RSQ_	_CP_	_PRESS_	Intercept
250	prog enz liver age gend alc1 alc2	7	0.22306	0.82129	14.2632	3.30254	4.41771
251	blood prog enz liver gend alc1 alc2	7	0.21081	0.84039	8.6793	2.94350	3.91149
252	blood prog enz liver age alc1 alc2	7	0.21136	0.83956	8.9214	2.90722	4.06637
253	blood prog enz liver age gend alc2	7	0.20867	0.84361	7.7352	2.88266	4.07191
254	blood prog enz age gend alc1 alc2	7	0.20705	0.84603	7.0295	2.80871	4.03678

----- Number of regressors in model=8 -----

Obs	VarsInModel	_IN_	_RMSE_	_RSQ_	_CP_	_PRESS_	Intercept
255	blood prog enz liver age gend alc1 alc2	8	0.20927	0.84613	9	2.93123	4.05052



# Best Models Using PRESS

- blood prog enz alc2, 2.73777
- blood prog enz age alc2, 2.73893
- blood prog enz age alc1 alc2, 2.77826
- blood prog enz age gend alc2, 2.77233
- blood prog enz gend alc2 2.78271
- blood prog enz alc1 alc2, 2.79565

# Automatic search procedures

- When  $p$  large, can't do all subset
- Stepwise type procedures
  - Forward selection (Step up)
  - Backward elimination (Step down)
  - Stepwise: combines forward and backward procedure allowing for removal of variables after adding new variables.
- Many other alternatives, but NONE guarantees the optimal solution (NP-hard problem)

# Forward Selection

- Start with an intercept
- At each step add the “best” variable (using some criteria –  $R^2$ , adj  $R^2$ ,  $C_p$ , partial correlation, SSE, ...)
- Compare the p-value for the test whether the just added variable is 0 with some pre-selected value.
  - If smaller – add the variable and repeat the procedure
  - If larger – stop. You have arrived at the final model.
- This is purely exploratory – see the book for comments

# Backwards Elimination

- Start with the full model
- At each step delete the “worst” variable (using some criteria – smallest increase in SSE, largest p-value,...)
- Compare the p-value for the test whether the just deleted variable is 0 with some pre-selected value.
  - If larger – delete the variable and repeat the procedure
  - If smaller – stop. You have arrived at the final model.

# Stepwise Regression

- Combines the forward and backward algorithm
- Starts at some model (empty or full is usual)
- First eliminates as many parameters as possible using backwards rule (using a minimum allowable p-value=  $\alpha$ -delete)
- Then attempt to add one variable (using a maximum allowable p-value= $\alpha$ -add)
- If a variable is added repeat, if not stop ( $\alpha$ -add <  $\alpha$ -delete required for convergence).

# Do it in SAS

```
* Forward selection;
proc reg data=surg;
    model lsurv=blood prog enz liver
    age gend alc1 alc2 /
    selection=FORWARD slentry=0.5;
run;
```

```
* backward elimination;
proc reg data=surg;
    model lsurv=blood prog enz liver
    age gend alc1 alc2 /selection=B
    slstay=0.05;
run;
```

```
* Forward stepwise;
proc reg data=surg;
    model lsurv=blood prog enz liver
    age gend alc1 alc2 /
    selection=STEPWISE slentry=0.50
    slstay=0.05;
run;
```

```
* Forward stepwise with first
    variable always in the model;
proc reg data=surg;
    model lsurv=blood prog enz liver
    age gend alc1 alc2 /
    selection=STEPWISE include=1;
run;
```

```
* Stepwise with at least two
    variables in the model;
proc reg data=surg;
    model lsurv=blood prog enz liver
    age gend alc1 alc2 /
    selection=STEPWISE start=2;
run;
```

# Forward Selection and Backward Elimination

---

No other variable met the 0.5000 significance level for entry into the model.

## Summary of Forward Selection

Step	Variable Entered	Number Vars In	Partial R-Square	Model R-Square	C(p)	F Value	Pr > F
1	enz	1	0.4276	0.4276	117.409	38.84	<.0001
2	prog	2	0.2357	0.6633	50.4716	35.70	<.0001
3	alc2	3	0.1147	0.7780	18.9145	25.85	<.0001
4	blood	4	0.0519	0.8299	5.7508	14.93	0.0003
5	gend	5	0.0076	0.8374	5.5406	2.23	0.1418
6	age	6	0.0060	0.8434	5.7874	1.80	0.1862
7	alc1	7	0.0026	0.8460	7.0295	0.77	0.3835

---

All variables left in the model are significant at the 0.0500 level.

## Summary of Backward Elimination

Step	Variable Removed	Number Vars In	Partial R-Square	Model R-Square	C(p)	F Value	Pr > F
1	liver	7	0.0001	0.8460	7.0295	0.03	0.8645
2	alc1	6	0.0026	0.8434	5.7874	0.77	0.3835
3	age	5	0.0060	0.8374	5.5406	1.80	0.1862
4	gend	4	0.0076	0.8299	5.7508	2.23	0.1418

# Stepwise

---

All variables left in the model are significant at the 0.0500 level.

The stepwise method terminated because the next variable to be entered was just removed.  
Summary of Stepwise Selection

Step	Variable Entered	Variable Removed	Number Vars In	Partial R-Square	Model R-Square	C(p)	F Value	Pr > F
1	enz		1	0.4276	0.4276	117.409	38.84	<.0001
2	prog		2	0.2357	0.6633	50.4716	35.70	<.0001
3	alc2		3	0.1147	0.7780	18.9145	25.85	<.0001
4	blood		4	0.0519	0.8299	5.7508	14.93	0.0003
5	gend		5	0.0076	0.8374	5.5406	2.23	0.1418
6		gend	4	0.0076	0.8299	5.7508	2.23	0.1418

---

All variables left in the model are required or significant at the 0.1500 level.

No other variable met the 0.1500 significance level for entry into the model.  
Summary of Stepwise Selection

Step	Variable Entered	Variable Removed	Number Vars In	Partial R-Square	Model R-Square	C(p)	F Value	Pr > F
1	enz		2	0.4424	0.6633	50.4716	67.01	<.0001
2	alc2		3	0.1147	0.7780	18.9145	25.85	<.0001
3	blood		4	0.0519	0.8299	5.7508	14.93	0.0003
4	gend		5	0.0076	0.8374	5.5406	2.23	0.1418



# Summary

- No method is the best for all model selection problems
- Consider more than one criterion
- “Best model” from automatic search procedures should be used as the starting point
- Apply knowledge of the subject matter to make a final selection – use your head!

# Model validation

- Three approaches to checking the validity of the model
  - Collect new data, does it fit the model
  - Compare with theory, other data, simulation
  - Use some of the data for the basic analysis and some for validity check, compare SSE with PRESS, MSE with MSPE