4 Introduction to Multiple Linear Regression

In **multiple linear regression**, a linear combination of two or more predictor variables is used to explain the variation in a response. In essence, the additional predictors are used to explain the variation in the response not explained by a simple linear regression fit.

As an illustration, I will consider the following problem. The data set is from the statistics package **Minitab**, where it is described thus:

Anthropologists wanted to determine the long-term effects of altitude change on human blood pressure. They measured the blood pressures of a number of Peruvians native to the high Andes mountains who had since migrated to lower climes. Previous research suggested that migration of this kind might cause higher blood pressure at first, but over time blood pressure would decrease. The subjects were all males over 21, born at high altitudes, with parents born at high altitudes. The measurements included a number of characteristics to help measure obesity: skin-fold and other physical characteristics. Systolic and diastolic blood pressure are recorded separately; systolic is often a more sensitive indicator. Note that this is only a portion of the data collected.

The data set is on the web site. Variables in the data set are

| Name | Description |
|---------|---------------------------|
| Age | Age in years |
| Years | Years since migration |
| Weight | Weight in kilograms |
| Height | Height in mm |
| Chin | Chin skin fold in mm |
| Forearm | Forearm skin fold in mm |
| Calf | Calf skin fold in mm |
| Pulse | Pulse in beats per minute |
| Systol | Systolic blood pressure |
| Diastol | Diastolic blood pressure |

A question we consider concerns the long term effects of an environmental change on the systolic blood pressure. In particular, is there a relationship between the systolic blood pressure and how long the Indians lived in their new environment as measured by the fraction of their life spent in the new environment? (fraction = years since migration/age - you need to generate fraction).

A plot of systolic blood pressure against fraction suggests a weak linear relationship (from graph matrix weight systol fraction). Nonetheless, consider fitting the regression model

sys
$$bp = \beta_0 + \beta_1 fraction + \epsilon$$
.

The least squares line is given by

$$sys \ bp = 133.50 - 15.75 \ fraction,$$

and suggests that average systolic blood pressure decreases as the fraction of life spent in modern society increases. However, the t-test of H_0 : $\beta_1 = 0$ is not significant at the 5% level (p-value=.089). That is, the weak linear relationship observed in the data is not atypical of a population where there is no linear relationship between systolic blood pressure and the fraction of life spent in a modern society.



Stata output:

| • | regress systol fraction | | | | | | | | | | | |
|---|-------------------------|--------------------------|----------------|------------|--------------------|----------------|--------------------------|----|--------------------|--|--|--|
| | Source | SS | df | | MS | | Number of obs $F(1)$ 37) | = | 39 3 05 | | | |
| | Model Residual | 498.063981 6033.37192 | 1 37 | 498 163 | .063981 .064106 | | Prob > F R-squared | = | 0.0888 | | | |
| | Total | 6531.4359 | 38 | 171 | .879892 | | Root MSE | = | 12.77 | | | |
| | systol | Coef. | Std. | Err. | t | P> t | [95% Conf. | In | terval] | | | |
| | fraction _cons | -15.75183 133.4957 | 9.012 4.038 | 962 011 | -1.75 33.06 | 0.089 0.000 | -34.01382 125.3139 | 2 | .510169 41.6775 | | | |

Even if this test were significant, the small value of $R^2 = .076$ suggests that fraction does not explain a **substantial** amount of the variation in the systolic blood pressures. If we omit the individual with the highest blood pressure (see the plot) then the relationship would be weaker.

Taking Weight into Consideration

At best, there is a weak relationship between systolic blood pressure and fraction. However, it is usually accepted that systolic blood pressure and weight are related; see the scatterplot matrix for confirmation. A natural way to take weight into consideration is to include weight and fraction as predictors of systolic blood pressure in the multiple regression model:

sys
$$bp = \beta_0 + \beta_1$$
 fraction $+\beta_2$ weight $+\epsilon$.

As in simple linear regression, the model is written in the form:

Response = Mean of Response + Residual,

so the model implies that average systolic blood pressure is a linear combination of fraction and weight. As in simple linear regression, the standard multiple regression analysis assumes that the responses are normally distributed with a constant variance $\sigma_{Y|X}^2$. The parameters of the regression model β_0 , β_1 , β_2 and $\sigma_{Y|X}^2$ are estimated by LS.

Stata output for fitting the multiple regression model follows.

| Source Model Residual | SS 3090.07324 3441.36266 | df 2 1545 36 95.5 | MS 5.03662 5934072 | | Number of obs F(2, 36) Prob > F R-squared Adj R-squared | $\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$ |
|-----------------------------|-----------------------------------|----------------------------------|--------------------------|-------------------------|--|--|
| Total | 6531.4359 | 38 171. | .879892 | | Root MSE | = 9.7772 |
| systol | Coef. | Std. Err. | t | P> t | [95% Conf. | Interval] |
| fraction weight _cons | -26.76722 1.216857 60.89592 | 7.217801 .2336873 14.28088 | -3.71 5.21 4.26 | 0.001 0.000 0.000 | -41.40559 .7429168 31.93295 | -12.12884 1.690796 89.85889 |

. regress systol fraction weight

Important Points to Notice About the Regression Output

1. The LS estimates of the intercept and the regression coefficient for fraction, and their standard errors, change from the simple linear model to the multiple regression model. For the simple linear regression

$$sys \ bp = 133.50 - 15.75 \ fraction.$$

For the multiple regression model

$$sys \ bp = 60.89 - 26.76 \ fraction + 1.21 \ weight.$$

There is frequently a big difference between coefficients from simple linear regression and those from multiple linear regression (for the *same* predictor variables).

- 2. Comparing the simple linear regression and the multiple regression models we see that the Model (Regression) df has increased to 2 from 1 (2=number of predictor variables) and the Residual (error) df has decreased from 37 to 36 (= n 1- number of predictors). Adding a predictor *increases* the Model (Regression) df by 1 and *decreases* the Residual df by 1.
- 3. The Residual SS decreases by 6033.37 3441.36 = 2592.01 upon adding the weight term. The Model (Regression) SS increased by 2592.01 upon adding the weight term to the model. The Total SS does not depend on the number of predictors so it stays the same. The Residual SS, or the part of the variation in the response unexplained by the regression model never increases when new predictors are added. After all, you are not going to do any *worse* modelling the data if you use more predictors the smaller model (simple linear regression) is a special case of the larger model (multiple linear regression). Anything you can fit using the simple one-variable model you also can fit using the two-variable model, but you can do a lot more with the two-variable model.
- 4. The proportion of variation in the response explained by the regression model:

 $R^2 = Model$ (or Regression) SS / Total SS

never decreases when new predictors are added to a model. The R^2 for the simple linear regression was .076, whereas $R^2 = .473$ for the multiple regression model. Adding the weight variable to the model increases R^2 by 40%. That is, weight and fraction together explain 40% more of the variation in systolic blood pressure than explained by fraction alone. I am not showing you the output, but if you predict systolic blood pressure using only weight, the R^2 is .27; adding fraction to *that* model increases the R^2 once again to .47. How well two predictors work together is not predictable from how well each works alone.

Stata also reports an *adjusted* R^2 . That has a penalty for fitting too many variables built into it, and can decrease when variables are added. If the number of variables is a lot less than n (it should be) there is not much difference between the two R^2 s.

5. The estimated variability about the regression line

Residual MS =
$$s_{Y|X}^2$$

decreased dramatically after adding the weight effect. For the simple linear regression model (fitting fraction as the only predictor), $s_{Y|X}^2 = 163.06$, whereas $s_{Y|X}^2 = 95.59$ for the multiple regression model. This suggests that an important predictor has been added to model. Note that *Stata* also reports Root MSE = $\sqrt{ResidualMS} = \sqrt{s_{Y|X}^2}$, an estimate of the standard deviation rather than the variance about the regression line.

6. The F-statistic for the multiple regression model

$$F_{obs} = \text{Regression MS} / \text{Residual MS} = 16.16$$

(which is compared to a F-table with 2 and 36 df) tests $H_0: \beta_1 = \beta_2 = 0$ against $H_A:$ not H_0 . This is a test of no relationship between the average systolic blood pressure and fraction and weight, assuming the relationship is linear. If this test is significant then either fraction or weight, or both, are important for explaining the variation in systolic blood pressure. Unlike simple linear regression, this test statistic is not simply the square of a *t*-statistic. It is a whole new test for us, and simply addresses the question "is anything going on anywhere in this model?"

7. Given the model

sys $bp = \beta_0 + \beta_1 fraction + \beta_2 weight + \epsilon$,

one interest is testing $H_0: \beta_2 = 0$ against $H_A: \beta_2 \neq 0$. The *t*-statistic for this test

$$t_{obs} = \frac{b_2 - 0}{SE(b_2)} = \frac{1.217}{.234} = 5.21$$

is compared to a t-critical value with Residual df = 36. Stata gives a p-value of .000, which suggests $\beta_2 \neq 0$. The *t*-test of H_0 : $\beta_2 = 0$ in the multiple regression model tests whether adding weight to the simple linear regression model explains a significant part of the variation in systolic blood pressure not explained by fraction. In some sense, the *t*-test of $H_0: \beta_2 = 0$ will be significant if the increase in R^2 (or decrease in Residual SS) obtained by adding weight to this simple linear regression model is substantial. We saw a big increase in R^2 , which is deemed significant by the *t*-test. A similar interpretation is given to the *t*-test for $H_0: \beta_1 = 0$.

- 8. The *t*-tests for $\beta_0 = 0$ and $\beta_1 = 0$ are conducted, assessed, and interpreted in the same manner. The p-value for testing $H_0: \beta_0 = 0$ is .000, whereas the p-value for testing $H_0: \beta_1 = 0$ is .001. This implies that fraction is important in explaining the variation in systolic blood pressure **after** weight is taken into consideration (by including weight in the model as a predictor). These *t*-tests are tests for the effect of a variable *adjusted* for the effects of all other variables in the model.
- 9. We compute CIs for the regression parameters β_i in the usual way: $b_i + t_{crit}SE(b_i)$, where t_{crit} is the *t*-critical value for the corresponding CI level with df = Residual df.

Understanding the Model

The *t*-test for $H_0: \beta_1 = 0$ is highly significant (p-value=.001) in the multiple regression model, which implies that fraction is important in explaining the variation in systolic blood pressure *after* weight is taken into consideration (by including weight in the model as a predictor). Weight is called a **suppressor variable**. Ignoring weight suppresses the relationship between systolic blood pressure and fraction - recall that fraction was not significant as a predictor by itself.

The implications of this analysis are enormous! Essentially, the correlation between a predictor and a response says very little about the importance of the predictor in a regression model with one or more additional predictors. This conclusion also holds in situations where the correlation is high, in the sense that a predictor that is highly correlated with the response may be unimportant in a multiple regression model once other predictors are included in the model.

Another issue that I wish to address concerns the interpretation of the regression coefficients in a multiple regression model. For our problem, let us first focus on the fraction coefficient in the fitted model

$$sys \ bp = 60.90 - 26.77 \ fraction + 1.22 \ weight.$$

The negative coefficient indicates that the predicted systolic blood pressure decreases as fraction increases **holding weight constant**. In particular, the predicted systolic blood pressure decreases by 26.76 for each unit increase in fraction, holding weight constant at any value. Similarly, the predicted systolic blood pressure increases by 1.21 for each unit increase in weight, holding fraction constant at any level.

We should examine residuals. Now the diagnostics are much more important to us, since we cannot see everything in terms of one predictor variable.



We will discuss these plots in class. Are there any observations we should investigate further? Which ones?

Another Multiple Regression Example

The data below are selected from a larger collection of data referring to candidates for the General Certificate of Education (GCE) who were being considered for a special award. Here, total denotes the candidate's TOTAL mark, out of 1000, in the GCE exam, while comp is the candidate's score in the compulsory part of the exam, which has a maximum score of 200 of the 1000 points on the exam. scel denotes the candidates' score, out of 100, in a School Certificate English Language (SCEL) paper taken on a previous occasion.

. list, clean

| | total | comp | scel |
|-----|-------|------|------|
| 1. | 476 | 111 | 68 |
| 2. | 457 | 92 | 46 |
| 3. | 540 | 90 | 50 |
| 4. | 551 | 107 | 59 |
| 5. | 575 | 98 | 50 |
| 6. | 698 | 150 | 66 |
| 7. | 545 | 118 | 54 |
| 8. | 574 | 110 | 51 |
| 9. | 645 | 117 | 59 |
| 10. | 690 | 114 | 80 |
| 11. | 634 | 130 | 57 |
| 12. | 637 | 118 | 51 |
| 13. | 390 | 91 | 44 |
| 14. | 562 | 118 | 61 |
| 15. | 560 | 109 | 66 |

A goal here is to compute a multiple regression of the TOTAL score on COMP and SCEL, and make the necessary tests to enable you to comment intelligently on the extent to which current performance in the compulsory test (COMP) may be used to predict aggregate TOTAL performance on the GCE exam, and on whether previous performance in the School Certificate English Language (SCEL) has any predictive value independently of what has already emerged from the current performance in the compulsory papers.

I will lead you through a number of steps to help you answer this question. Let us answer the following straightforward questions based on the **Stata** output.

1. Plot TOTAL against COMP and SCEL individually, and comment on the form (i.e. linear, non-linear, logarithmic, etc.), strength, and direction of the relationships.

2. Plot COMP against SCEL and comment on the form, strength, and direction of the relationship.

3. Compute the correlation between all pairs of variables. Do the correlation values appear reasonable, given the plots?

Stata output: scatterplot matrix and correlations...



In parts 4 through 9, ignore the possibility that TOTAL, COMP or SCEL might ideally need to be transformed.

0.0346

0.0527

4. Which of COMP and SCEL explains a larger proportion of the variation in TOTAL? Which would appear to be a better predictor of TOTAL? (Explain).

5. Consider 2 simple linear regression models for predicting TOTAL one with COMP as a predictor, and the other with SCEL as the predictor. Do COMP and SCEL individually appear to be important for explaining the variation in TOTAL (i.e. test that the slopes of the regression lines are zero). Which, if any, of the output, support, or contradicts, your answer to the previous question?

Stata output:

| Source | | df | | MS | | Number of obs | = | 15 |
|-------------------|--------------------------|----------------|---------------|---|----------------|--------------------------|-------------|--------------------|
| Model Residual | 53969.7272 47103.2062 | 1 13 | 53969 3623 | 9.7272 .32355 | | Prob > F R-squared | = | 0.0020 |
| Total | 101072.933 | 14 | 7219 | . 49524 | | Root MSE | = | 60.194 |
| total | Coef. | Std. | Err. | t | P> t | [95% Conf. | In | terval] |
| comp _cons | 3.948465 128.5479 | 1.023 115.1 | 073 604 | 3.86 1.12 | 0.002 0.285 | 1.73825 -120.241 | 6 3 | .158681 77.3367 |
| . regress tota | al scel | | | | | | | |
| Source | SS SS | df | | MS | | Number of obs $F(1)$ 13) | = | 15 |
| Model Residual | 30320.6397 70752.2936 | 1 13 | 30320 5442 |).6397 .48412 | | Prob > F R-squared | - - - | 0.0346 |
| Total | 101072.933 | 14 | 7219 | . 49524 | | Root MSE | = | 73.773 |
| total | Coef. | Std. | Err. | t | P> t | [95% Conf. | In | terval] |
| scel _cons | 4.826232 291.5859 | 2.044 119.0 | 738 382 | $\begin{array}{c} 2.36\\ 2.45\end{array}$ | 0.035 0.029 | $.4088448 \\ 34.4196$ | 9 5 | .243619 48.7522 |
| | | | | | | | | |

. regress total comp

6. Fit the multiple regression model

$$\Gamma \text{OTAL} = \beta_0 + \beta_1 \text{COMP} + \beta_2 \text{SCEL} + \epsilon.$$

Test $H_0: \beta_1 = \beta_2 = 0$ at the 5% level. Describe in words what this test is doing, and what the results mean here.

Stata output:

. regress total comp scel

| Source | SS | df | MS | | Number of obs | = 15 = 8 1/ |
|---------------------------|---------------------------------|----------------------------------|----------------------|-------------------------|------------------------------------|----------------------------------|
| Model Residual | 58187.5043 42885.429 | 2 2909 12 3573 | 3.7522 .78575 | | Prob > F R-squared | = 0.0058 = 0.5757 = 0.5050 |
| Total | 101072.933 | 14 7219 | .49524 | | Root MSE | = 59.781 |
| total | Coef. | Std. Err. | t | P> t | [95% Conf. | Interval] |
| comp scel _cons | 3.295936 2.09104 81.16147 | 1.180318 1.924796 122.4059 | 2.79 1.09 0.66 | 0.016 0.299 0.520 | .7242444 -2.102731 -185.5382 | 5.867628 6.284811 347.8611 |

7. In the multiple regression model, test $H_0: \beta_1 = 0$ and $H_0: \beta_2 = 0$ individually. Describe in words what these tests are doing, and what the results mean here.

8. How does the R^2 from the multiple regression model compare to the R^2 from the individual simple linear regressions? Is what you are seeing here appear reasonable, given the tests on the individual coefficients?

9. Do your best to answer the question posed above, in the paragraph on page 43 that begins "A goal". Provide an equation (LS) for predicting TOTAL.

Comments on the GCE Analysis

I will give you my thoughts on these data, and how I would attack this problem, keeping the ultimate goal in mind. As a first step, I plot the data and check whether transformations are needed. The plot of TOTAL against COMP is fairly linear, but the trend in the plot of TOTAL against SCEL is less clear. You might see a non-linear trend here, but the relationship is not very strong. When I assess plots I try to not allow a few observations affect my perception of trend, and with this in mind, I do not see any strong evidence at this point to transform any of the variables.

One difficulty that we must face when building a multiple regression model is that these twodimensional (2D) plots of a response against individual predictors may have little information about the appropriate scales for a multiple regression analysis. In particular, the 2D plots only tell us whether we need to transform the data in a simple linear regression analysis. If a 2D plot shows a strong non-linear trend, I would do an analysis using the suggested transformations, including any other effects that are important. However, it might be that no variables need to be transformed in the multiple regression model.

Although SCEL appears to be useful as a predictor of TOTAL on its own, the multiple regression output indicates that SCEL does not explain a significant amount of the variation in TOTAL, once the effect of COMP has been taken into account. In particular, the SCEL effect in the multiple regression model is far from significant (p-value=.30). Hence, previous performance in the SCEL exam has little predictive value independently of what has already emerged from the current performance in the compulsory papers.

What are my conclusions? Given that SCEL is not a useful predictor in the multiple regression model, I would propose a simple linear regression model to predict TOTAL from COMP:

Predicted TOTAL = 128.55 + 3.95 COMP.

Output from the fitted model was given earlier. A residual analysis of the model showed no serious deficiencies. In particular, the residuals versus the predicted values looks random and the normal probability plot of the residuals looks reasonably straight. Note that the following summaries are for this one-variable model, not the two-variable model fit earlier.

| | Shap | oiro-Wilk W | test for | normal data | |
|----------|-------|-------------|----------|-------------|---------|
| Variable | l Obs | W | V | Z | Prob>z |
| residual | 15 | 0.97287 | 0.526 | -1.271 | 0.89806 |





A Taste of Model Selection for Multiple Regression

Given data on a response variable Y and k predictor variables $X_1, X_2, ..., X_k$, we wish to develop a regression model to predict Y. Assuming that the collection of variables is measured on the correct scale, and that the candidate list of predictors includes all the important predictors, the most general (linear) model is

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \epsilon.$$

In most problems one or more of the predictors can be eliminated from this general or **full model** without loss of information. We want to identify the important predictors, or equivalently, eliminate the predictors that are not useful for explaining the variation in Y.

We will study several automated methods for model selection. Given a specific criterion for selecting a model, **Stata** gives the best predictors. Before applying any of the methods, you should plot Y against each predictor $X_1, X_2, ..., X_k$ to see whether transformations are needed. If a transformation of X_i is suggested, include the transformation along with the original X_i in the candidate list. Note that you can transform the predictors differently, for example, $log(X_1)$ and $\sqrt{X_2}$. However, if several transformations are suggested for the response, then you should consider doing one analysis for each suggested response scale before deciding on the final scale.

At this point, I will only consider the **backward elimination method**. Other approaches can be handled in **Stata**.

Backward Elimination

The backward elimination procedure deletes unimportant variables, one at a time, starting from the full model. The steps in the procedure are:

1. Fit the full model

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \epsilon.$$
(1)

2. Find the variable which when omitted from the full model (1) reduces R^2 the least, or equivalently, increases the Residual SS the least. This is the variable that gives the largest p-value for testing an individual regression coefficient $H_0: \beta_i = 0$ for i > 0. Suppose this variable is X_k . If you reject H_0 , stop and conclude that the full model is best. If you do not reject H_0 , delete X_k from the full model, giving the new full model

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_{k-1} X_{k-1} + \epsilon.$$

Repeat steps 1 and 2 sequentially until no further predictors can be deleted.

In backward elimination we isolate the least important predictor left in the model, and check whether it is important. If not, delete it and repeat the process. Otherwise, stop. A test level of 0.1 (a very common value to use), for example, on the individual predictors is specified in **Stata** using pr(0.1) in the sw command.

Epidemiologists use a slightly different approach to building models. They argue strongly for the need to always include **confounding** variables in a model, regardless of their statistical significance. I will discuss this issue more completely for logistic regression, but you should recognize its importance.

Illustration

I will illustrate backward elimination on the Peru Indian data, using systolic blood pressure as the response, and seven candidate predictors: weight in kilos, height in mm, chin skin fold in mm, forearm skin fold in mm, calf skin fold in mm, pulse rate-beats/min, and fraction. A plot of systolic blood pressure against each of the individual potential predictors does not strongly suggest the need to transform either the response or any of the predictors:

| | | • | • | | • |
|---------------------------------------|---------------------------------------|--------------|------|----------------|-----------------|
| Weight | · | · | | h: : ::::::::: | 169. |
| *** · | Height | | | | |
| | Chin | | | | 4 3. • • • • |
| su ³ | ~~~ | Forearm | | 1 | · |
| | | . <u>B</u> . | Calf | | |
| · · · · · · · · · · · · · · · · · · · | · · · · · · · · · · · · · · · · · · · | • | P | ulse | · · · · · · |
| | | | | fraction | |
| | | • | | . . | Systol |

The correlation matrix shows that most of the potential predictors are weakly correlated with systolic blood pressure. Based on correlations, the best single variable for predicting blood pressure is weight.

. correlate weight height chin forearm calf pulse fraction systol (obs=39)

| | weight | height | chin | forearm | calf | pulse | fraction | systol |
|--|--|---|--------------------------------------|----------------------------|-------------------|------------------|-------------------|--------|
| weight height chin forearm calf pulse | 1.0000 0.4503 0.5617 0.5437 0.3919 0.3118 | 1.0000 -0.0079 -0.0689 -0.0028 0.0078 | 1.0000 0.6379 0.5160 0.2231 | 1.0000 0.7355 0.4219 | 1.0000 0.2087 | 1.0000 | | |
| fraction systol | 0.2931 0.5214 | 0.0512 0.2191 | $0.1201 \\ 0.1702$ | 0.0280 0.2723 | -0.1130 0.2508 | 0.2142 0.1355 | 1.0000 -0.2761 | 1.0000 |

Stata commands for the previous output are (assuming you grabbed the Stata data set peru.dta from the web and use'd it)

generate fraction=years/age graph matrix weight height chin forearm calf pulse fraction systol correlate weight height chin forearm calf pulse fraction systol

Summaries from the full model with 7 predictors follow. The *F*-test in the full model ANOVA table (F = 4.92 with p-value=.0008) tests the hypothesis that the regression coefficient for each predictor variable is zero. This test is highly significant, indicating that one or more of the predictors is important in the model. Note that $R^2 = .53$ for the full model.

| · TEBLEDD DAD | or werging mer | gine chimite | Jiearm car | T harse | TIACTION | |
|---|---|---|---|--|--|---|
| Source | SS | df | MS | | Number of obs | = 39 |
| Model Residual | 3436.89993 3094.53596 | 7 490. 31 99.8 | 985705 3237407 | | Prob > F R-squared Adi B-squared | = 4.92 = 0.0008 = 0.5262 = 0.4192 |
| Total | 6531.4359 | 38 171. | 879892 | | Root MSE | = 9.9912 |
| systol | Coef. | Std. Err. | t | P> t | [95% Conf. | Interval] |
| weight height chin forearm calf pulse fraction _cons | $\begin{array}{c} 1.710538 \\0454089 \\ -1.154889 \\7143249 \\ .1058654 \\ .07971 \\ -29.35489 \\ 106.3085 \end{array}$ | .3864434 .0394397 .845932 1.350676 .6116778 .1959149 7.86754 53.8376 | $\begin{array}{r} 4.43 \\ -1.15 \\ -1.37 \\ -0.53 \\ 0.17 \\ 0.41 \\ -3.73 \\ 1.97 \end{array}$ | 0.000 0.258 0.182 0.601 0.864 0.687 0.001 0.057 | .9223814 1258466 -2.880179 -3.469047 -1.14166 3198611 -45.40084 -3.494025 | 2.498694 .0350289 .5704004 2.040397 1.35339 .4792811 -13.30894 216.111 |

. regress systol weight height chin forearm calf pulse fraction

You can automate stepwise selection of predictors (for which backward elimination is a special case) using the **sw** command. Six model selection procedures are allowed: backward selection, forward selection, backward stepwise, forward stepwise, backward hierarchical selection, and forward hierarchical selection. See the **Stata** manual for descriptions. The command **sw** can also be used with other regression models including logistic (and other binary response model) regression, Poisson regression, and Cox proportional hazards regression. To obtain the stepwise procedure for multiple linear regression in our example, using an cutoff of 0.1, type **sw regress systol weight** height chin forearm calf pulse fraction, pr(0.1). I cannot seem to get this to work correctly using the pull-down menus, and I'm not sure there is much potential gain anyway. This is a pretty simple command. The **Stata** output follows.

| . sw regress p = 0.8637 >= p = 0.6953 >= p = 0.6670 >= p = 0.2745 >= p = 0.1534 >= | systol weight begi = 0.1000 remo = 0.1000 remo = 0.1000 remo = 0.1000 remo = 0.1000 remo | height c n with fu ving calf ving puls ving fore ving heig ving chin | hin forear ll model e arm ht | m calf p | pulse fraction, | pr(0.1) |
|---|--|--|--|-------------------------|-----------------------------------|-----------------------------------|
| Source | SS | df | MS | | Number of obs | = 39 |
| Model Residual | 3090.07324 3441.36266 | 2 15 36 95 | 45.03662 .5934072 | | Prob > F R-squared | = 0.0000 = 0.4731 |
| Total | 6531.4359 | 38 17 | 1.879892 | | Root MSE | = 9.7772 |
| systol | Coef. | Std. Err | . t | P> t | [95% Conf. | [Interval] |
| weight fraction _cons | 1.216857 -26.76722 60.89592 | .2336873 7.217801 14.28088 | 5.21 -3.71 4.26 | 0.000 0.001 0.000 | .7429168 -41.40559 31.93295 | 1.690796 -12.12884 89.85889 |

The procedure summary tells you that the least important variable in the full model, as judged by the p-value, is calf skin fold. This variable, upon omission, will reduce R^2 the least, or equivalently, increases the Residual SS the least. The p-value of .86 exceeds the specified 0.10 cut-off, so the first step of the backward elimination would be to eliminate calf skin fold from the model. This is the p-value for the t-test on calf in the 7-variable model.

The next variable eliminated is pulse because of the p-value of .70 in the 6-variable model where calf was not fit (**Stata** isn't showing you all of that output). Notice that this is different from the p-value in the 7-variable model. Next **Stata** removes forearm because of the large p-value of .67 in a 5-variable model with calf and pulse removed. Other variables are eliminated similarly. There is a huge amount of computation summarized in this one table.

Looking at the rest of the step history, the backward elimination procedure eliminates five variables from the full model, in the following order: calf skin fold, pulse rate, forearm skin fold, height, and chin skin fold. As we progress from the full model to the selected model, R^2 decreases as follows: $R^2 = .53$ (full model), .53, .52, .52, .50, and .47 (from several regression fits not shown). The decrease is slight across this spectrum of models.

The model summary selected by backward elimination includes two predictors: weight and fraction. The fitted model is given by:

Predicted SYS BP = 60.90 + 1.22 Weight - 26.77 Fraction.

Each predictor is significant at the .001 level. The fitted model explains 47% of the variation in systolic blood pressures. This 2-variable model does as well, for any practical purposes, in predicting systolic blood pressure as a much more complicated 7-variable model. There was no real surprise here, since these two variables were the only ones significant in the 7-variable model, but often you will be left with a model you would not have guessed from a fit of all variables.

Using a mechanical approach, we are led to a model with weight and fraction as predictors of systolic blood pressure. At this point you should closely examine the fitted model.

Stepwise procedures receive a great deal of criticism. When a large number of variables are screened this way, the resulting relationships tend to be exaggerated. There is a big multiple comparisons problem here as well. This technique should be regarded as exploratory and the resulting p-values and coefficients assessed from independent data, although common practice is just to report final results. It is likely that the strength of relationships discovered in stepwise procedures will be hard to replicate in later studies, however. This is, nonetheless, an invaluable screening device when one has a lot of predictor variables.