

SPSS FAQ

How can I compare regression coefficients between two groups?

Sometimes your research hypothesis may stipulate that the size of a regression coefficient should be bigger for one group than for another. For example, you might believe that the regression coefficient of height predicting weight would be higher for men than for women. Below, we have a data file with 10 fictional females and 10 fictional males, along with their height in inches and their weight in pounds.

```
data list free
/ id * gender (AB) height * weight.
begin data.
1 F 56 117
2 F 60 125
3 F 64 133
4 F 68 141
5 F 72 149
6 F 74 109
7 F 62 128
8 F 65 131
9 F 65 131
10 F 70 145
11 M 64 211
12 M 68 223
13 M 72 235
14 M 76 247
15 M 80 259
16 M 82 262
17 M 69 228
18 M 74 241
19 M 75 241
20 M 82 269
end data.
```

We analyzed their data separately using the regression commands below. Note that we have to do two regressions, one with the data for females only and one with the data for males only. We can use the `split` file command to split the data file by gender and then run the regression. The parameter estimates (coefficients) for females and males are shown below, and the results do seem to suggest that height is a stronger predictor of weight for males (3.16) than for females (2.00).

```
sort cases by gender.
split file by gender.
regression
/dep weight
/method = enter height.
split file off.
```

Variables Entered/Removed ^a				
gender	Model	Variables Entered	Variables Removed	Method
F	1	height		Enter
M		height		Enter

a. All requested variables entered.
b. Dependent Variable: weight

Model Summary					
gender	Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
F	1	.884 ^a	.780	.776	9.31141
M	1	.884 ^a	.808	.807	2.40738

a. Predictors: (Constant), height

ANOVA ^a							
gender	Model		Sum of Squares	df	Mean Square	F	Sig.
F	1	Regression	1318.561	1	1318.561	359.812	.000 ^b
		Residual	29.339	8	3.667		
		Total	1348.900	9			
M	1	Regression	3082.536	1	3082.536	669.926	.000 ^b
		Residual	46.364	8	5.795		
		Total	3128.900	9			

a. Predictors: (Constant), height
b. Dependent Variable: weight

Coefficients ^a									
gender	Model		Unstandardized Coefficients			Standardized Coefficients			Sig.
			B	Std. Error	t	Beta	1		
F	1	(Constant)	-2.387	7.053			.340	.743	
		height	2.006	1.10	.889		.18985	.000	
M	1	(Constant)	5.852	8.930			.627	.548	
		height	3.160	1.23	.884		.25383	.000	

a. Dependent Variable: weight

We can compare the regression coefficients of males with females to test the null hypothesis $H_0: \beta_{0F} = \beta_{0M}$, where β_{0F} is the regression coefficient for females, and β_{0M} is the regression coefficient for males. To do this analysis, we first make a dummy variable called female that is coded 1 for female and 0 for male, and a variable female that is the product of female and height. We then use female, height and female*height as predictors in the regression equation.

```
split file off.
compute female = 0.
if gender = "F" female = 1.
compute femht = female*height.
execute.
regression
/dep weight
/method = enter female height femht.
```

The output is shown below:

Variables Entered/Removed ^a				
Model	Variables Entered	Variables Removed	Method	
1	female, height, femht		Enter	

a. All requested variables entered.
b. Dependent Variable: weight

Model Summary					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	
1	.884 ^a	.809	.808	2.11818	

a. Predictors: (Constant), femht, height, female

ANOVA ^a						
Model	Sum of Squares	df	Mean Square	F	Sig.	
1	Regression	3032.705	3	1010.902	4250.111	.000 ^b
	Residual	75.762	6	12.627		
	Total	3108.467	9			

a. Predictors: (Constant), femht, height, female
b. Dependent Variable: weight

Coefficients ^a									
Model			Unstandardized Coefficients			Standardized Coefficients			Sig.
			B	Std. Error	t	Beta	1		
1	(Constant)		5.852	8.930			.694	.491	
	female		-7.899	11.371	-.703		-.703	.492	
	height		3.160	1.11	4.21		.28446	.000	
	femht		-.084	.168	-.508		-.620	.000	

a. Dependent Variable: weight

The term femht tests the null hypothesis $H_0: \beta_{0F} = \beta_{0M}$. The T value is -5.02 and is significant, indicating that the regression coefficient β_{0F} is significantly different from β_{0M} .

Let's look at the parameter estimates to get a better understanding of what they mean and how they are interpreted. First, recall that our dummy variable female is 1 if female and 0 if male; therefore, males are the omitted group. This is needed for proper interpretation of the estimates.

Parameter	
Variable	Estimate
INTERCEPT 5.60147 - This is the intercept for the males (omitted group)	
This corresponds to the intercept for males in the separate groups analysis.	
FEMALE -7.899161 - Intercept Female - Intercept males	
This corresponds to difference of the intercepts from the separate groups analysis, and is indeed -2.39747046 - 5.60147149	
HEIGHT 3.169721 - Slope for males (omitted group), i.e., β_{0F} .	
FEMHT -1.093895 - Slope for female - Slope for males (i.e., $\beta_{0F} - \beta_{0M}$)	
From the separate groups, this is indeed 2.06872170 - 3.16972143 .	

It is also possible to run such an analysis using `glm`, using syntax like that below. Note that other statistical packages, such as SAS and Stata, omit the group of the dummy variable that is coded as zero. However, SPSS omits the group coded as one. Therefore, when you compare the output from the different packages, the results seem to be different. To make the SPSS results match those from other packages, you need to create a new variable that has the opposite coding (i.e., switching the zeros and ones). We do this with the male variable. We do not know of an option in SPSS `glm` to easily change which group is the omitted group. Please note that you can use the `contrast` subcommand to get the contrast coefficient for female using 2 as the reference group; however, the coding of female in the interaction is such that 1 is used as the reference group, so the use of the `contrast` subcommand is not very helpful in this situation.

```
compute male = not female.
glm weight by male with height
/design = male weight male by height
/print = parameter.
```

Between-Subjects Factors		
male	count	total
1.00	10	

Tests of Between-Subjects Effects						
Dependent Variable: weight	Source	Type III Sum of Squares	df	Mean Square	F	Sig.
	Corrected Model	3032.705	3	1010.902	4250.111	.000
	Intercept	374	1	374	.074	.782
	male	1342	1	1342	.490	.482
	height	495.831	1	495.831	.992	.000
	male*height	201.115	1	201.115	.420	.000
	Error	75.762	6	12.627		
	Total	73114.000	20			
	Corrected Total	6345.892	19			

a. R Squared = .999 (Adjusted R Squared = .999)

Parameter Estimates						
Dependent Variable: weight	Parameter	B	Std. Error	t	Sig.	95% Confidence Interval
	Intercept	5.602	8.969	.624	.537	[-11.624, 22.797]
	male=0.00	-7.899	11.371	-.703	.482	[-32.094, 16.895]
	male=1.00	.0 ^a				
	height	3.169	1.11	28.446	.000	[2.994, 3.424]
	male=0.00*height	-1.094	.168	-6.520	.000	[-1.450, -.738]
	male=1.00*height	.0 ^a				

a. This parameter is set to zero because it is redundant.

SPSS FAQ

How can I compare regression coefficients across three (or more) groups?

Sometimes your research hypothesis may predict that the size of a regression coefficient may vary across groups. For example, you might believe that the regression coefficient of height predicting weight would differ across three age groups (young, middle age, senior citizens). Below, we have a data file with 10 fictional young people, 10 fictional middle age people, and 10 fictional senior citizens, along with their height in inches and their weight in pounds. The variable `age` indicates the age group and is coded 1 for young people, 2 for middle aged, and 3 for senior citizens. Below we show two ways that you can get this data file into SPSS. One way is to cut and paste the following code into an SPSS syntax window and run it.

```
data list list / id age height weight.
```

```
begin data.
```

```
1 1 56 140
```

```
2 1 60 155
```

```
3 1 64 143
```

```
4 1 68 161
```

```
5 1 72 139
```

```
6 1 54 159
```

```
7 1 62 138
```

```
8 1 65 121
```

```
9 1 65 161
```

```
10 1 70 145
```

```
11 2 56 117
```

```
12 2 60 125
```

```
13 2 64 133
```

```
14 2 68 141
```

```
15 2 72 149
```

```
16 2 54 109
```

```
17 2 62 128
```

```
18 2 65 131
```

```
19 2 65 131
```

```
20 2 70 145
```

```
21 3 64 211
```

```
22 3 68 223
```

```
23 3 72 235
```

```
24 3 76 247
```

```
25 3 80 259
```

```
26 3 62 201
```

```
27 3 69 228
```

```
28 3 74 245
```

```
29 3 75 241
```

```
30 3 82 269
```

```
end data.
```

```
execute.
```

Another way is to click on `compress.sav` and then use the `get file` command (insert the proper drive letter if you did not place the file in your current directory).

```
get file 'c:\compress.sav'.
```

After first sorting by age, we analyze the data for each age group separately using the `regression` command. In order to use just the data for a specific age group, we need to use a filter to "filter out" the other data. Remember that when you have completed the analysis, you need to turn the filter off.

```
sort cases by age.
```

```
split file by age.
```

```
regression
```

```
  /dep weight
```

```
  /method=enter height.
```

```
split file off.
```

```
exe.
```

The parameter estimates (coefficients) for the young, middle age, and senior citizens are shown below, and the results do seem to suggest that height is a stronger predictor of weight for seniors (3.18) than for the middle aged (2.09). The results also seem to suggest that height does not predict weight as strongly for the young (-.37) as for the middle aged and seniors. However, we would need to perform specific significance tests to be able to make claims about the differences among these regression coefficients.

< some output omitted to save space >

Model Summary					
age	Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
100	1	.712 ^a	.509	-.393	11.42000
200	1	.889 ^a	.878	.876	1.91504
300	1	.984 ^a	.988	.987	2.40138

a. Predictors: (Constant), height

ANOVA ^b						
age	Model		Sum of Squares	df	Mean Square	F
100	1	Regression	42.857	1	42.857	237
		Residual	1440.143	9	160.110	849 ^c
		Total	1483.000	9		
200	1	Regression	1319.561	1	1319.561	358.812
		Residual	26.339	9	2.867	.000 ^d
		Total	1345.900	9		
300	1	Regression	3882.536	1	3882.536	699.928
		Residual	45.364	9	5.045	.000 ^d
		Total	3927.900	9		

a. Predictors: (Constant), height

b. Dependent Variable: weight

Coefficients ^a					
age	Model		Unstandardized Coefficients	Standardized Coefficients	
100	1	(Constant)	170.166	49.430	3.443
		height	-.377	.714	-.487
200	1	(Constant)	-2.307	7.063	-.340
		height	2.026	1.10	.889
300	1	(Constant)	5.602	8.558	.627
		height	3.150	1.23	.984

a. Dependent Variable: weight

We can compare the regression coefficients among these three age groups to test the null hypothesis

H0: $\beta_1 = \beta_2 = \beta_3$

where β_1 is the regression for the young, β_2 is the regression for the middle aged, and β_3 is the regression for senior citizens. To do this analysis, we first make a dummy variable called `age1` that is coded 1 if young (`age=1`), 0 otherwise, and `age2` that is coded 1 if middle aged (`age=2`), 0 otherwise. We also create `age1ht` that is `age1` times height, and `age2ht` that is `age2` times height.

```
compute age1 = 0.
```

```
compute age2 = 0.
```

```
if age = 1 age1 = 1.
```

```
if age = 2 age2 = 1.
```

```
compute age1ht = age1*height.
```

```
compute age2ht = age2*height.
```

```
execute.
```

We can now use `age1`, `age2`, `age1ht`, and `age2ht` as predictors in the regression equation in the `regress` command below. The `regression` command will be followed by

```
/method = test(age1 age2)
```

```
and
```

```
/method = test(age1ht age2ht)
```

The first one provides a 2 degree of freedom test to determine if, taken together, the variable `age` is statistically significant. We have included this for the sake of completeness, because this is a standard part of the analysis. The second subcommand tests the null hypothesis

H0: $\beta_1 = \beta_2 = \beta_3$

This test will also have 2 degrees of freedom because it compares among three regression coefficients.

```
regression
  /dep weight
  /method = enter height
  /method=test(age1 age2)
  /method = test(age1ht age2ht).
```

< some output omitted to save space >

Model Summary					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	
1	.714 ^a	.510	.492	25.2749	
2	.974 ^a	.940	.942	11.83281	
3	.989 ^a	.979	.974	7.84844	

a. Predictors: (Constant), height

b. Predictors: (Constant), height, age2, age1

c. Predictors: (Constant), height, age2, age1, age2ht, age1ht

ANOVA ^b							
Model			Sum of Squares	df	Mean Square	F	Sig.
1		Regression	3626.967	1	3626.967	26.141	.000 ^a
		Residual	3496.033	28	124.859		
		Total	7112.000	29			
2	Subset	age1, age2	3143.844	2	1571.922	109.360	.000 ^a
			3740.811	3	2246.937	157.803	.000 ^a
	Total	Regression	3702.189	28	132.200		
		Residual	7112.000	29			
		Total	7112.000	29			
3	Subset	age1ht, age2ht	2106.544	2	1053.272	17.292	.000 ^a
			2106.544	2	1053.272	17.292	.000 ^a
	Total	Regression	3996.395	5	1399.071	220.261	.000 ^a
		Residual	1515.605	24	63.150		
		Total	7112.000	29			

a. Predictors: (Constant), height

b. Tested against the full model

c. Predictors in the Full Model: (Constant), height, age2, age1

d. Predictors in the Full Model: (Constant), height, age2, age1, age2ht, age1ht

e. Dependent Variable: weight

Coefficients ^a					
Model			Unstandardized Coefficients	Standardized Coefficients	
1	(Constant)		158.151	85.362	-2.403
		height	6.959	.714	.510
2	(Constant)		108.492	27.745	3.910
		height	1.195	.301	.445
	age1		-74.524	6.261	-.722
		age2	-89.524	6.261	-.876
3	(Constant)		5.602	28.488	1.90
		height	3.190	.407	.459
	age1		184.595	41.555	1.593
		age2	-7.989	41.555	-.077
	age1ht		-3.567	.813	-.208
		age2ht	-1.024	.813	-.077

a. Dependent Variable: weight

The analysis below shows that the null hypothesis

H0: $\beta_1 = \beta_2 = \beta_3$

SPSS FAQ

How do I interpret the parameter estimates for dummy variables in regression or glm?

Consider this simple data file that has nine subjects (sub) in three groups (iv) with a score on the outcome or dependent variable (dv).

```
data list list / sub iv dv.
begin data
1 1 48
2 1 49
3 1 50
4 2 17
5 2 20
6 2 23
7 3 28
8 3 30
9 3 32
end data.
```

Below we use the **means** command to find the overall mean and the means for the three groups.

```
means tables = dv by iv.
```

As we see below, the overall mean is 33, and the means for groups 1, 2 and 3 are 49, 20 and 30 respectively.

Case Processing Summary					
DV * IV	Cases		Excluded		Total
	N	Percent	N	Percent	N
DV * IV	9	100.0%	0	.0%	9
Report					
IV	Mean	N	Std. Deviation		
1.00	49.0000	3	1.00000		
2.00	20.0000	3	3.00000		
3.00	30.0000	3	2.00000		
Total	33.0000	9	12.89380		

Let's run a standard ANOVA on these data using glm.

```
glm dv by iv.
```

The results of the ANOVA are shown below.

Between-Subjects Factors					
IV	1.00		N		3
	2.00				3
	3.00				3
Tests of Between-Subjects Effects					
Dependent Variable: DV					
Source	Type III Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	1302.000(a)	2	651.000	139.500	.000
Intercept	9801.000	1	9801.000	2100.214	.000
IV	1302.000	2	651.000	139.500	.000
Error	28.000	6	4.667		
Total	11131.000	9			
Corrected Total	1330.000	8			
a. R Squared = .979 (Adjusted R Squared = .972)					

Now, let's take this information we have found and relate it to the results that we get when we run a similar analysis using dummy coding. Let's make a data file called **dummy2** that has dummy variables called **iv1** (1 if iv=1), **iv2** (1 if iv=2) and **iv3** (1 if iv=3). Note that **iv3** is not really necessary, but it could be useful for further exploring the meaning of dummy variables. We will then use the **regression** command to predict **dv** from **iv1** and **iv2**.

```
compute iv1 = 0.
if iv = 1 iv1 = 1.
compute iv2 = 0.
if iv = 2 iv2 = 1.
compute iv3 = 0.
if iv = 3 iv3 = 1.
execute.

regression
/dependent = dv
/method = enter iv1 iv2.
```

The output is shown below.

Variables Entered/Removed(b)						
Model	Variables Entered		Variables Removed	Method		
1	IV2, IV1(a)			Enter		
a. All requested variables entered.						
b. Dependent Variable: DV						
Model Summary						
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate		
1	.989(a)	.979	.972	2.16525		
a. Predictors: (Constant), IV2, IV1						
ANOVA(b)						
Model	Sum of Squares		df	Mean Square	F	Sig.
1	Regression	1302.000	2	651.000	139.500	.000(a)
	Residual	28.000	6	4.667		
	Total	1330.000	8			
a. Predictors: (Constant), IV2, IV1						
b. Dependent Variable: DV						
Coefficients(a)						
Model	Unstandardized Coefficients			Standardized Coefficients		
	B	Std. Error	Beta	t	Sig.	
1	(Constant)	30.000	1.247		24.054	
	IV1	19.000	1.764		737	10.772
	IV2	-10.000	1.764		-388	-5.669
a. Dependent Variable: DV						

First, note that from the ANOVA using the **glm** command that the **F**-value was 139.5 and for the regression using the **regression** command the **F**-value (for the model) is also 139.5. This illustrates that the overall test of the model using regression is really the same as doing an ANOVA.

After the **ANOVA** table, there is a table entitled **Coefficients**. What is the interpretation of the values listed there, the 30, 19 and -10? Notice how we have **iv1** and **iv2** that refer to group 1 and group 2, but we did not include any dummy variable referring to group 3. Group 3 is often called the **omitted group** or **reference group**. Recall that the means of the 3 groups were 49, 20 and 30 respectively. The **Intercept** term is the mean of the dependent variable, which we called **dv**, for the **omitted group**, and indeed the parameter estimate (in the column **B**) from the output is the mean of group 3, 30. The parameter estimate for **iv1** is the mean of the dependent variable, **dv**, for group 1 minus the mean of the dependent variable for group 3, 49 - 30 = 19, and indeed that is the parameter estimate for **iv1**. Likewise, the parameter estimate for **iv2** is the mean of the dependent variable for group 2 minus the mean of the dependent variable for group 3, 20 - 30 = -10, the parameter estimate for **iv2**.

So, in summary:

Intercept	mean of group 3 (mean of omitted group)
iv1	mean of group 1 - group 3 (omitted group)
iv2	mean of group 2 - group 3 (omitted group)

Try running this example, but use **iv2** and **iv3** using **regression** (making group 1 the omitted group) and see what happens.

Finally, consider how the parameter estimates can be used in the regression model to obtain the means for the groups (the predicted values).

The regression model is

$$\hat{y}_{\text{predicted}} = 30 + iv1 \cdot 19 + iv2 \cdot -10$$

For group 1: $\hat{y}_{\text{predicted}} = 30 + 1 \cdot 19 + 0 \cdot -10 = 49$
For group 2: $\hat{y}_{\text{predicted}} = 30 + 0 \cdot 19 + 1 \cdot -10 = 20$
For group 3: $\hat{y}_{\text{predicted}} = 30 + 0 \cdot 19 + 0 \cdot -10 = 30$

As you see, the regression formula predicts that each group will have the mean value of its group.

You can also perform the same analysis using **glm**. The **print = parameter** subcommand tells SPSS to print the regression coefficients.

```
glm dv with iv1 iv2
/print = parameter.
```

Tests of Between-Subjects Effects					
Dependent Variable: DV					
Source	Type III Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	1302.000(a)	2	651.000	139.500	.000
Intercept	2700.000	1	2700.000	578.571	.000
IV1	941.500	1	941.500	116.036	.000
IV2	150.000	1	150.000	32.143	.001
Error	28.000	6	4.667		
Total	11131.000	9			
Corrected Total	1330.000	8			
a. R Squared = .979 (Adjusted R Squared = .972)					

Parameter Estimates					
Dependent Variable: DV					
Parameter	B	Std. Error	t	Sig.	95% Confidence Interval
Intercept	30.000	1.247	24.054	.000	26.948 33.052
IV1	19.000	1.764	10.772	.000	14.684 23.316
IV2	-10.000	1.764	-5.669	.001	-14.316 -5.684

SPSS FAQ

How can I perform hypothesis tests in glm?

Sometimes you may want to test hypotheses about the parameters after a linear regression analysis. On this page, we show a couple of examples of how to perform these hypothesis tests using the **lmatrix** and **kmatrix** subcommands in the **glm** procedure. These examples will use data set [hsb2.sav](#). Let's say that we have run a linear regression model as follows:

```
glm write with female read math
/print=parameter
/design=female read math.
```

Parameter Estimates

Dependent Variable: write score						
Parameter df	B	Std. Error	t	Sig.	95% Confidence Interval Lower Bound	Upper Bound
Intercept	11.896	2.883	4.155	.000	6.250	17.542
female	5.443	.935	5.822	.000	3.599	7.287
read	.325	.081	5.355	.000	.205	.445
math	.397	.086	5.986	.000	.267	.528

Written as a regression equation, we have the following:

write = **b_0** + **b_1*** **female** + **b_2*****read** + **b_3*****math**, where **b_0** = 11.896, **b_1** = 5.443, **b_2** = .325 and **b_3** = .397.

Example 1

Let's say that we want to test if the coefficient for **read** is equal to the coefficient for **math**. The **lmatrix** subcommand allows us to specify our hypothesis test in terms of the linear combination of the regression coefficients. In our case, our null hypothesis is that **b_2** = **b_3**, or equivalently, **b_2**-**b_3** = 0. This leads to our **lmatrix** subcommand with 1 following the variable **read** and -1 following the variable **math**.

```
glm write with female read math
/print=parameter
/design=female read math
/lmatrix = 'math = read' read 1 math -1.
```

Custom Hypothesis Tests

Contrast Results (K Matrix)^a

Contrast				Dependen... writing score
L1	Contrast Estimate			-.072
	Hypothesized Value			0
	Difference (Estimate - Hypothesized)			-.072
	Std. Error			.116
	Sig.			.534
	95% Confidence Interval for Difference			-.301
	Lower Bound			Upper Bound .156

a. Based on the user-specified contrast coefficients (L1) matrix: math = read

Test Results

Dependent Variable: writing score					
Source	Sum of Squares	df	Mean Square	F	Sig.
Contrast	16.792	1	16.792	.388	.534
Error	8473.526	196	43.232		

In the output, we see the difference between the two parameters is $-.072 = (.325 - .397)$, as we expected. What the output also gives is the standard error for the difference and the confidence interval. The Test Results table shows the F-value and the p-value.

Example 2

Let's say that we want to test if the coefficient for **female** is equal to 4.2. In order to do this, we need to use the **kmatrix** subcommand, because we are testing if the value is something other than 0. You might want to do this, if, for example, you had regression coefficients from a previous model and you wanted to see if they were equal to the coefficients obtained with your current model. To keep the example simple, we will test only one variable (**female**) in this example.

```
glm write with female read math
/print=parameter
/design=female read math
/lmatrix = 'female' female 1
/kmatrix 4.2.
```

Contrast Results (K Matrix)^a

Contrast				Dependen... writing score
L1	Contrast Estimate			5.443
	Hypothesized Value			4.200
	Difference (Estimate - Hypothesized)			1.243
	Std. Error			.935
	Sig.			.185
	95% Confidence Interval for Difference			-.601
	Lower Bound			Upper Bound 3.087

a. Based on the user-specified contrast coefficients (L1) matrix: female

Test Results

Dependent Variable: writing score					
Source	Sum of Squares	df	Mean Square	F	Sig.
Contrast	76.452	1	76.452	1.768	.185
Error	8473.526	196	43.232		

Example 3

Let's say that we want to test if the coefficient for **female** is equal to 4.2 and that the coefficient for **read** is equal to the coefficient for **math**. This will be a two degree-of-freedom test since there are two hypotheses that we want to test simultaneously. Notice that the values specified on the **kmatrix** subcommand are listed in the same order as the tests listed on the **lmatrix** subcommand.

```
glm write with female read math
/print=parameter
/design=female read math
/lmatrix = 'test' female 1; read 1 math -1
/kmatrix 4.2; 0.
```

Contrast Results (K Matrix)^a

Contrast				Dependen... writing score
L1	Contrast Estimate			5.443
	Hypothesized Value			4.200
	Difference (Estimate - Hypothesized)			1.243
	Std. Error			.935
	Sig.			.185
	95% Confidence Interval for Difference			-.601
	Lower Bound			Upper Bound 3.087
L2	Contrast Estimate			-.072
	Hypothesized Value			.000
	Difference (Estimate - Hypothesized)			-.072
	Std. Error			.116
	Sig.			.534
	95% Confidence Interval for Difference			-.301
	Lower Bound			Upper Bound .156

a. Based on the user-specified contrast coefficients (L1) matrix: test

Test Results

Dependent Variable: writing score					
Source	Sum of Squares	df	Mean Square	F	Sig.
Contrast	95.324	2	47.662	1.102	.334
Error	8473.526	196	43.232		

SPSS FAQ

How can I analyze multiple mediators in SPSS?

This page was created by Christopher J. Preacher of the University of Kansas and Andrew F. Hayes of Ohio State University. The code can be found on the sidebar accompanying their paper:

Preacher, K. J., & Hayes, A. F. (2008). *Asymptotic and resampling strategies for assessing and comparing indirect effects in multiple mediator models*. *Behavior Research Methods*, 40, 878-891.

To use the Preacher and Hayes' script, first download the file here. The file for the documentation describing how to download and install the .sbs file is here. Then open it in SPSS and run it by clicking on the green arrow or choosing "Run" from the Macro menu. This will open a SPSS dialog window.

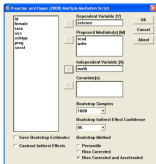


Example 1: Multiple mediators

For this example, we will use the `hsb2` dataset with `science` as the dependent variable, `math` as the independent variable and `read` and `write` as the two mediator variables. The paths in such a model are depicted below. In our analysis, we are interested in finding these paths to calculate the direct and indirect effects of our variables.



To begin, we indicate which of our variables are the dependent, independent, and mediator variables in the dialog window.



This generates the output below.

```
Run MATRIX procedure:
Dependent, Independent, and Proposed Mediator Variables:
IV = math
DV = science
MEDI = read
      write

Sample size
200

IV to Mediators (a paths)
      Coeff      se      t      p
read      .1048      .0087    12.4178    .0000
write     .0566      .0066    11.0452    .0000

Direct Effects of Mediators on DV (b paths)
      Coeff      se      t      p
read      .0587      .0087    6.7112    .0000
write     .0393      .0078    5.0385    .0000

Total Effect of IV on DV (c path)
      Coeff      se      t      p
math      .0645      .0083    11.4371    .0000

Direct Effect of IV on DV (c-prime path)
      Coeff      se      t      p
math      .0318      .0077    4.1654    .0000

Model Summary for DV Model
R-sq Adj R-sq      F      df1      df2      p
.4999      .4923    65.3187    3,0000    196.0000    .0000

=====
NORMAL THEORY TESTS FOR INDIRECT EFFECTS

Indirect Effects of IV on DV through Proposed Mediators (ab paths)
Effect      se      z      p
TOTAL      .3476      .0596    5.8277    .0000
read       .2186      .0524    4.1692    .0000
write      .1290      .0454    2.8422    .0045

=====
BOOTSTRAP RESULTS FOR INDIRECT EFFECTS

Indirect Effects of IV on DV through Proposed Mediators (ab paths)
      Effect      Boot      Bias      SE
TOTAL      .3476      .3449    -.0027    .0645
read       .2186      .2164    -.0022    .0537
write      .1290      .1293     .0003    .0496

Bias Corrected and Accelerated Confidence Intervals
Lower      Upper
TOTAL      .2230      .4700
read       .1125      .3243
write      .0294      .2235

=====
Level of Confidence for Confidence Intervals:
95

Number of Bootstrapping Resamples:
1000

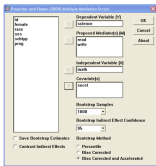
----- END MATRIX -----
```

The results above assuming normality suggest that each of the separate indirect effects as well as the total indirect effect are significant. From the above results it is also possible to compute the ratio of indirect to direct effect (.3476/.3187 = 1.09) and the proportion of the total effect due to the indirect effect (.3476/.6451 = 53.9% = .54).

The normal theory tests for indirect effects compute the standard errors using the delta method which assumes that the estimates of the indirect effects are normally distributed. For many situations this is acceptable, but does not work well for the indirect effects which are usually positively skewed and variable. Thus the bias and accelerated CI methods are more appropriate for these situations. However, it is recommended that the normal distribution errors and confidence intervals be used. Additionally, if your outcome is binary, a logit or a probit, bootstrap estimates should be used. These can be found in the next block of output. These standard errors are slightly larger than those indicated assuming normality and the normal interpretation remains the same.

Example 2: Multiple mediators with control variable(s)

What do you do if you also have control variables? You add them as covariates to the model. Let's say that `score2` is a covariate. We can see the new results below.



```
Run MATRIX procedure:
Dependent, Independent, and Proposed Mediator Variables:
IV = math
DV = science
MEDI = read
      write

Statistical Controls:
score2, score3

Sample size
200

IV to Mediators (a paths)
      Coeff      se      t      p
read      .0538      .0634    0.9512    .0000
write     .1144      .0617    1.7112    .0000

Direct Effects of Mediators on DV (b paths)
      Coeff      se      t      p
read      .1100      .0729    1.5117    .0000
write     .1119      .0748    1.4973    .0045

Total Effect of IV on DV (c path)
      Coeff      se      t      p
math      .0671      .0484    1.3863    .0000

Direct Effect of IV on DV (c-prime path)
      Coeff      se      t      p
math     -.0218      .0773   -0.1654    .0000

Partial Effect of Control Variables on DV
      Coeff      se      t      p
score2    -.0227      .0645   -0.3516    .7251

Model Summary for DV Model
R-sq Adj R-sq      F      df1      df2      p
.3502      .3400    48.8010    4,0000    195.0000    .0000

=====
BOOTSTRAP RESULTS FOR INDIRECT EFFECTS

Indirect Effects of IV on DV through Proposed Mediators (ab paths)
      Effect      Boot      Bias      SE
TOTAL      .2452      .2439    -.0014    .0478
read       .1102      .1039    -.0063    .0413
write      .0891      .0900    -.0009    .0342

Bias Corrected and Accelerated Confidence Intervals
Lower      Upper
TOTAL      .1197      .3522
read       .0863      .2476
write      .0311      .1671

=====
Level of Confidence for Confidence Intervals:
95

Number of Bootstrapping Resamples:
1000

----- END MATRIX -----
```

SPSS FAQ

How can I get out-of-sample predicted values?

Sometimes it is useful to get predicted values for cases that were not used in the regression analysis. There are two ways to do this in SPSS. Let's use the `hsb2` dataset and create some missing values in a variable. Specifically, we will set the first nine values in the variable `write` to be missing. Then we will use `write` as our outcome variable in an OLS regression analysis. Of course, the cases with missing values will not be used in the analysis, but we can still get the predicted values for those cases.

```
get file = 'd:/data/hsb2.sav'.
```

```
sort cases by id.
if id lt 10 write = $sysmis.
list write read math
/cases=from 1 to 12.
```

write	read	math
.	34.00	40.00
.	39.00	33.00
.	63.00	48.00
.	44.00	41.00
.	47.00	43.00
.	47.00	46.00
.	57.00	59.00
.	39.00	52.00
.	48.00	52.00
54.00	47.00	49.00
46.00	34.00	45.00
44.00	37.00	45.00

Number of cases read: 12 Number of cases listed: 12

Method 1

When running the `regression` command, we can use the `save` subcommand to save the predicted values to the current data file. We have supplied the name for the new variable in parentheses after the SPSS keyword `pred`. After running the regression, we will list the first 12 cases in the data set for the variables `write` and `pred_1`.

```
regression
/dependent write
/method = enter read math
/save pred(pred_1).
```

<output omitted>

```
list write pred_1
/cases from 1 to 12.
```

write	pred_1
.	42.24554
.	40.81015
.	54.03857
.	45.58411
.	47.28941
.	48.53128
.	56.83733
.	48.67533
.	51.30748
54.00	49.77315
46.00	44.31532
44.00	45.19271

Number of cases read: 12 Number of cases listed: 12

Method 2

Another way to get out-of-sample predictions is to save the model information to an `.xml` file, use the `model handle` command to name the `.xml` file, and then use the `ApplyModel` function of the `compute` command to create the predicted values. We will list the first 12 cases in the data file for the variables `write` and `yhat`.

```
regression
/dependent write
/method = enter read math
/outfile=model('d:/data/working/hsb_ml.xml').
```

<output omitted>

```
model handle name = ml file='d:/data/working/hsb_ml.xml'.
```

```
compute yhat = ApplyModel(ml,'predict').
```

```
list write yhat
/cases from 1 to 12.
```

write	yhat
.	42.25
.	40.81
.	54.04
.	45.58
.	47.29
.	48.53
.	56.84
.	48.68
.	51.31
54.00	49.77
46.00	44.32
44.00	45.19

Number of cases read: 12 Number of cases listed: 12

Now let's look at `pred_1` and `yhat` side by side; as you can see, they are the same.

```
formats pred_1 yhat (f8.5).
```

```
list write pred_1 yhat
/cases from 1 to 12.
```

write	pred_1	yhat
.	42.24554	42.24554
.	40.81015	40.81015
.	54.03857	54.03857
.	45.58411	45.58411
.	47.28941	47.28941
.	48.53128	48.53128
.	56.83733	56.83733
.	48.67533	48.67533
.	51.30748	51.30748
54.00	49.77315	49.77315
46.00	44.31532	44.31532
44.00	45.19271	45.19271

Number of cases read: 12 Number of cases listed: 12

SPSS FAQ

How can I output the results of my regression to an SPSS data file?

Sometimes it is useful to output the results of a regression analysis to a data file for further analyses. To do this in SPSS, you can use the **output** subcommand of the **regression** command. You have two choices of what to save using this subcommand: you can save the covariance matrix of the coefficients (with the **covb** option) or you can save the correlation matrix of the coefficients (with the **corb** option). Let us use a data set called **hsb2** as an example. We will save the results and the covariance matrix of the coefficients in a file called **out1.sav**.

```
get file 'd:\hsb2.sav'.
regression
  /dep = write
  /method = enter read female
  /outfile = covb('d:\out1.sav').
```

Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	20.228	2.714		.000
	READ	.566	.049	.612	.000
	FEMALE	5.487	1.014	.289	.000

a. Dependent Variable: WRITE

```
get file 'd:\out1.sav'.
list.
```

DEPVAR_	ROWTYPE_	VARNAME_	CONST_	READ	FEMALE
WRITE	COV	CONST_	7.36	-.13	-.70
WRITE	COV	READ	-.13	.00	.00
WRITE	COV	FEMALE	-.70	.00	1.03
WRITE	EST		20.23	.57	5.49
WRITE	SE		2.71	.05	1.01
WRITE	SIG		.00	.00	.00
WRITE	DFE		197.00	197.00	197.00

Number of cases read: 7 Number of cases listed: 7

As you can see above, the covariances between the estimates have been saved to this file, as well as the estimates, their standard errors, the significance and the error degrees of freedom. Note that the precision of the values saved in **out1.sav** is greater than the two decimal places shown here. Two decimal places are shown because that is the default number of decimal places to display in an SPSS data set. You can easily increase the number of decimal places shown by going to the "Variable View" of the SPSS Data Editor and increasing the value in the column labeled "Decimals" (you may have to increase the column width first).

Now let's run the same regression and this time use the **corb** option instead.

```
get file 'd:\hsb2.sav'.
regression
  /dep = write
  /method = enter read female
  /outfile = corb('d:\out2.sav').
```

Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	20.228	2.714		.000
	READ	.566	.049	.612	.000
	FEMALE	5.487	1.014	.289	.000

a. Dependent Variable: WRITE

```
get file 'd:\out2.sav'.
list.
```

DEPVAR_	ROWTYPE_	VARNAME_	CONST_	READ	FEMALE
WRITE	COR	CONST_	1.00	-.96	-.25
WRITE	COR	READ	-.96	1.00	.05
WRITE	COR	FEMALE	-.25	.05	1.00
WRITE	EST		20.23	.57	5.49
WRITE	SE		2.71	.05	1.01
WRITE	SIG		.00	.00	.00
WRITE	DFE		197.00	197.00	197.00

Number of cases read: 7 Number of cases listed: 7

As we can see, the correlations between the coefficients have been saved to the data set, as well as the estimates, their standard errors, the significance and the error degrees of freedom.

Results can also be saved into data sets using the Output Management System (OMS). For more information on how to use OMS to output results to data sets, please see [How can I output my results to a data file in SPSS?](#).

SPSS FAQ

How can I test a group of variables in SPSS regression?

Suppose that you want to run a regression model and to test the statistical significance of a group of variables. For example, let's say that you want to predict students' writing score from their reading, math and science scores. The data set with these variables in it can be downloaded by following this link: [hds2.dat](#).

The SPSS syntax for this would be:

```
regression
  /dependent = write
  /method = enter read math science.
```

Variables Entered/Removed(b)				
Model	Variables Entered	Variables Removed	Method	
1	science score, reading score, math score(a)		. Enter	

a All requested variables entered.

b Dependent Variable: writing score

Model Summary				
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.684(a)	.467		459

a Predictors: (Constant), science score, reading score, math score

ANOVA(b)					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	8353.990	3	2784.663	57.302	.000(a)
Residual	9524.885	196	48.596		
Total	17878.875	199			

a Predictors: (Constant), science score, reading score, math score

b Dependent Variable: writing score

Coefficients(a)					
Model	Unstandardized Coefficients			Standardized Coefficients	
	B	Std. Error		Beta	t
(Constant)	13.192	3.069			4.299 .000
1 reading score	.236	.069		.255	3.410 .001
math score	.319	.076		.316	4.222 .000
science score	.202	.069		.211	2.918 .004

a Dependent Variable: writing score

Now let's suppose that you wanted to test the combined effect of math and science on writing. The SPSS syntax for doing that is below. Note that the variables listed in the **method = test()** subcommand are not listed on the **method = enter** subcommand. In other words, the independent variables are listed only once. Also note that, unlike other SPSS subcommands, you can have multiple **method =** subcommands within the **regression** command.

```
regression
  /dependent = write
  /method = enter read
  /method = test(math science).
```

Variables Entered/Removed(b)				
Model	Variables Entered	Variables Removed	Method	
1	reading score(a)		. Enter	
2	science score, math score		. Test	

a All requested variables entered.

b Dependent Variable: writing score

Model Summary				
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.597(a)	.356	.353	7.62487
2	.684(b)	.467	.459	6.97111

a Predictors: (Constant), reading score

b Predictors: (Constant), reading score, science score, math score

ANOVA(c)							
Model	Sum of Squares	df	Mean Square	F	Sig.	R Square Change	
1 Regression	6367.421	1	6367.421	109.521	.000(a)		
Residual	11511.454	198	58.139				
Total	17878.875	199					
Subset Tests	1986.569	2	993.284	20.439	.000(b)	.111	
2 Regression	8353.990	3	2784.663	57.302	.000(c)		
Residual	9524.885	196	48.596				
Total	17878.875	199					

a Predictors: (Constant), reading score

b Tested against the full model.

c Predictors in the Full Model: (Constant), reading score, science score, math score.

d Dependent Variable: writing score

Coefficients(a)					
Model	Unstandardized Coefficients			Standardized Coefficients	
	B	Std. Error		Beta	t
(Constant)	23.959	2.808			8.539 .000
1 reading score	.552	.053		.597	10.465 .000
(Constant)	13.192	3.069			4.299 .000
2 reading score	.236	.069		.255	3.410 .001
math score	.319	.076		.316	4.222 .000
science score	.202	.069		.211	2.918 .004

a Dependent Variable: writing score

Excluded Variables(b)					
Model	Beta In	t	Sig.	Partial Correlation	Collinearity Statistics
					Tolerance
1 math score	-.395(a)	5.583	.000	.370	.561
science score	-.322(a)	4.609	.000	.312	.663

a Predictors in the Model: (Constant), reading score

b Dependent Variable: writing score

If you wanted to test all three variables together, the syntax would be:

```
regression
  /dependent = write
  /method = test(read math science).
```

Variables Entered/Removed(a)				
Model	Variables Entered	Variables Removed	Method	
1	science score, reading score, math score		. Test	

a Dependent Variable: writing score

Model Summary				
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.684(a)	.467		459

a Predictors: (Constant), science score, reading score, math score

ANOVA(c)					
Model	Sum of Squares	df	Mean Square	F	Sig.
Subset Tests	reading score, math score, science score	8353.990	3	2784.663	57.302 .000(a)
1 Regression	8353.990	3	2784.663	57.302 .000(b)	
Residual	9524.885	196	48.596		
Total	17878.875	199			

a Tested against the full model.

b Predictors in the Full Model: (Constant), science score, reading score, math score.

c Dependent Variable: writing score

Coefficients(a)					
Model	Unstandardized Coefficients			Standardized Coefficients	
	B	Std. Error		Beta	t
(Constant)	13.192	3.069			4.299 .000
1 reading score	.236	.069		.255	3.410 .001
math score	.319	.076		.316	4.222 .000
science score	.202	.069		.211	2.918 .004

a Dependent Variable: writing score

You will notice that the output from the first example with the three independent variables on the **method = enter** subcommand and the output from this example with the three independent variables on the **method = test()** subcommand are virtually identical. The only difference between them is the line in the ANOVA table that gives the test of the subset, which in this case is all of the variables. The point of this example is that you can put all of the independent variables in the regression on the **method = test()** subcommand and not use a **method = enter** subcommand if you like.

SPSS FAQ

How can I create a scatterplot with a regression line in SPSS?

There are at least two ways to make a scatterplot with a regression line in SPSS. One way is to use the **graph** command, and another way is to use the **ggraph** command. Both are illustrated below.

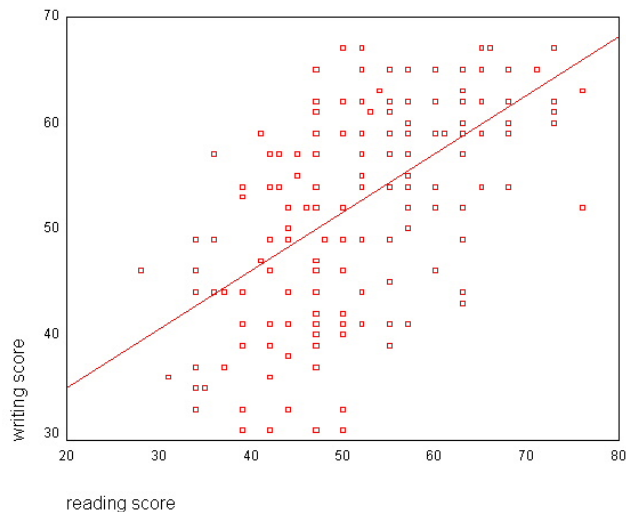
Let's read in an example dataset, [hsb2](#), which contains data from the High School and Beyond study.

```
get file 'c:\hsb2.sav'.
```

Let's do a scatterplot of the variables **write** with **read**. Note that after running the code below, you need to double click on the graph, which will open up the chart editor window. Select "chart" from the menu at the top and then "options" (which is the first item in the menu). On the right of the dialogue box is a check box called "Total" under the heading "Fit Line". Click on the box to put in the check, or click on "Fit Options" to select a different type of fit method, such as lowess, quadratic or cubic. Close the chart editor so that the changes take effect on your graph.

Using the graph command

```
graph  
  /scatterplot(bivar)=read with write.
```



Using the ggraph command

The **ggraph** command was introduced in version 14 of SPSS. This command can be used to create and edit scatterplots. Below is the syntax for creating a scatterplot with the regression line.

```
GGRAPH  
  /GRAPHDATASET NAME="graphdataset" VARIABLES=read write  
  /GRAPHSPEC SOURCE=INLINE.  
BEGIN GPL  
  SOURCE: s=userSource(id("graphdataset"))  
  DATA: read=col(source(s), name("read"))  
  DATA: write=col(source(s), name("write"))  
  GUIDE: axis(dim(1), label("reading score"))  
  GUIDE: axis(dim(2), label("writing score"))  
  ELEMENT: line( position( smooth.linear(read*write) ) )  
  ELEMENT: point(position(read*write))  
END GPL.
```

