

Bivariate Analysis Instructions

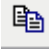
- Bivariate Analysis Instructions 1
- Bivariate Data Analysis using Linear Regression and Genstat 1
 - Predictions 3
 - Correlations 4
 - Finding the typical data- using Summary Stats 4
 - Piecewise Functions 5
 - Plotting more than one pair of variables at the same time -matrix 7
 - Linear Regression with groups 8
 - Residuals 9
 - Outliers..... 10
 - Restricting/Filtering Data 10
 - Non- Linear Models 10
 - Exponential Function 11
 - Power function 13
 - Polynomial 15

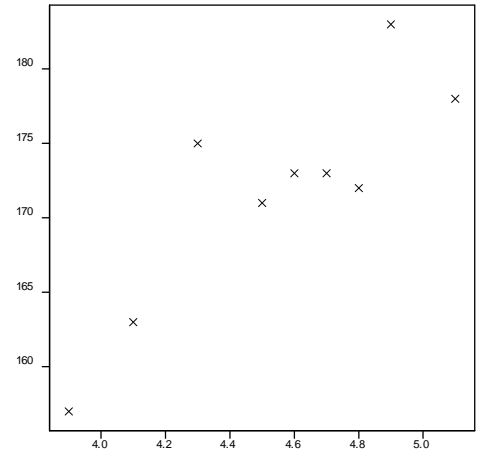
Bivariate Data Analysis using Linear Regression and Genstat




1. Open Genstat
2. Open the file *metacarpal*
3. You should get this menu
4. Just click on Finish and your file will be in Genstat
5. To draw a scatterplot of the data, use the pull-down **Graphics** menu and select **2D Scatter Plot**
6. Fill in as shown by double clicking on the variables and then clicking **Run**.

Row	metacarpal_bone_length_in_cm	stature_in_cm
1	4.5	171
2	5.1	178
3	3.9	157
4	4.1	163
5	4.8	172
6	4.9	183

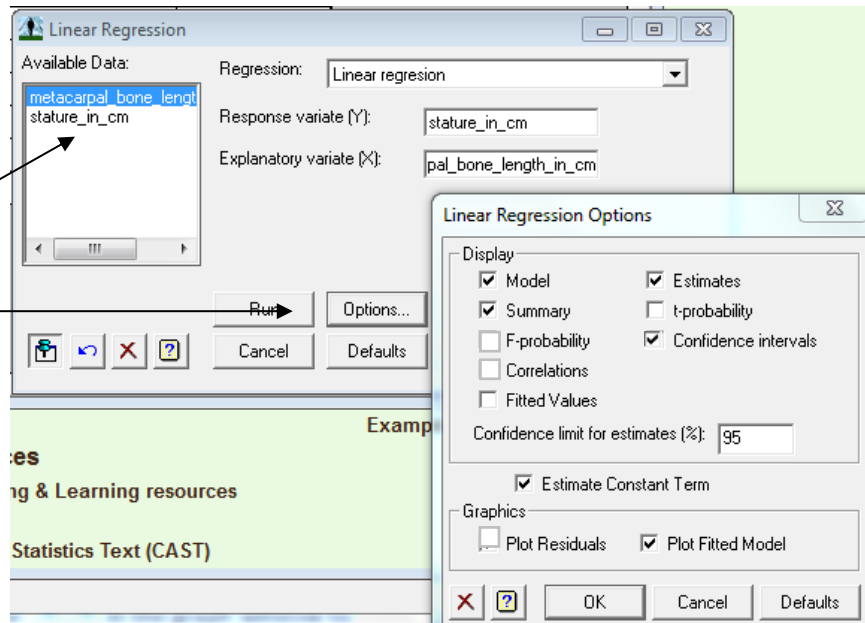
7. You should now get the graph! A **right click** will give the option to copy or click on  and the graph can be pasted into a Word document.



8. To return to the Spreadsheet, click on the  icon along the task bar at the bottom of the screen.

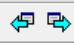
9. To perform the linear regression, use the **Stats** menu and select **Linear Regression**.


10. Fill in the dialogue box as shown, double clicking on the variables to select them. Click on **Options** to select further options and select by clicking. Fill in as shown.




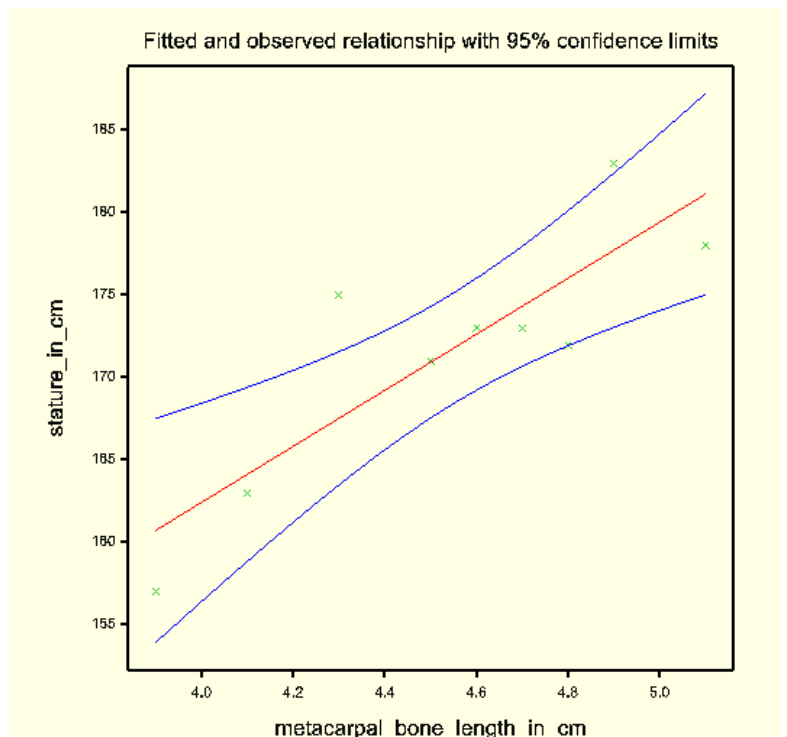
11. Click **OK** and then **Run**.

12. You will now get a graph of the fitted

model. Use the  in the graph window to move between graphs. To

find the output, click on the  and under the **Window** menu, select **Output**. This can be copied into Word, though you will need to select the regression output you require first.

To return to the graphs at any time just click on the 



Regression analysis

Response variate: stature_in_cm

Fitted terms: Constant, metacarpal_bone_l_length_in_cm

Message: the following units have high leverage.

Unit	Response	Leverage
3	157.00	0.46

this means that this point has a big effect on the trend line and hence the regression equation.

Estimates of parameters

Parameter	estimate	s.e.	t(7)	t pr.
Constant	94.4	17.7	5.34	0.001
metacarpal_bone_length_in_cm	17.00	3.88	4.38	0.003

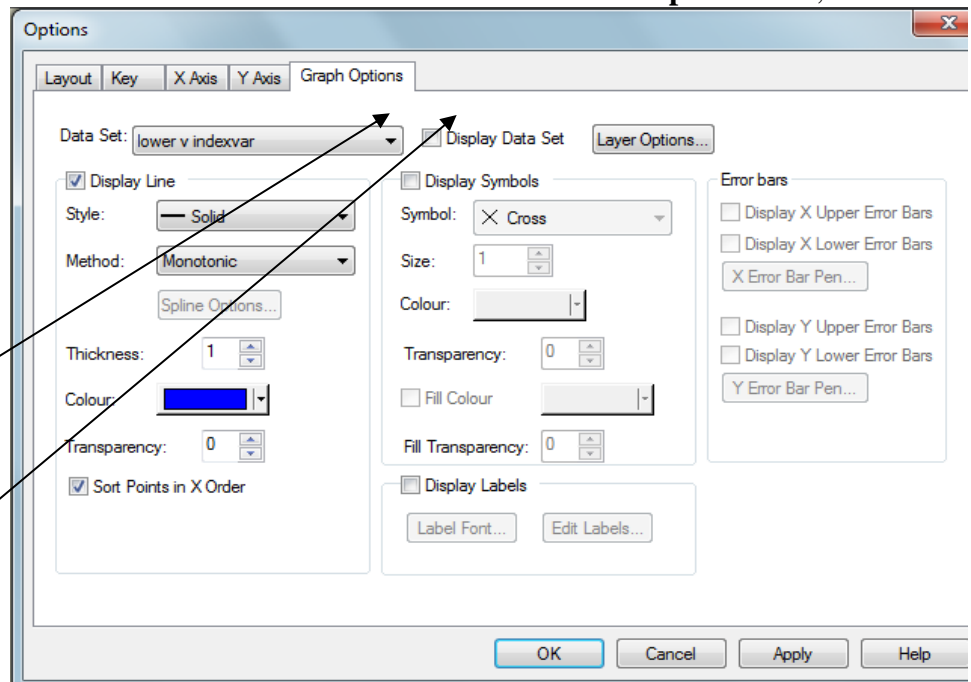
The model is stature = 17 x metacarpal bone length + 94.4 cm

The graphs can be edited to remove the confidence levels if desired. In the **Graph** window, chose **Edit** and then **Edit Graph**.

You now choose **Edit** and then **Graph Options**
By choosing the two Data set

- Lower v indexvar
- Upper v indexvar

and clicking off **Display data set** you remove the lines.



Predictions


You can use your model to predict the height when given

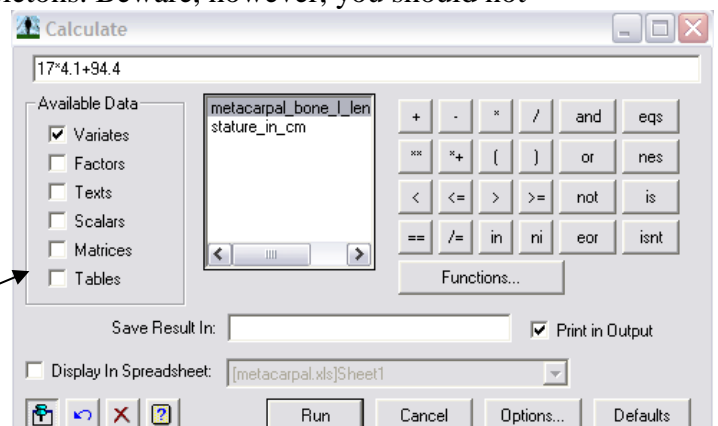
the length of the metacarpal bone for other skeletons. Beware, however, you should not

extrapolate (use your model outside the x values used to produce your model) only to

intrapolate – find a height for a bone length **inside** the bone length range you used to produce the model. For example find the predicted height for a bone length of 4.1

a. either work it out manually by substituting 4.1 into the equation,
 $height = 17 \times bone + 94.4$ or


b. use Genstat Calculator  by typing in as shown



you should get $(17 \times 4.1) + 94.4$ and selecting print in Output, you will get **164.1** in the Output

Correlations

To find r under the **Stats** menu choose **correlations** and then **correlation coefficient**

1. Click on  to put your variables in the Data column, tick on **Correlations** to ensure that you get the correlations
2. Click Run

Correlations between parameter estimates

Parameter	ref	correlations
Constant	1	1.000
metacarpal_bone_length_in_cm	2	-0.997 1.000
	1	2

Genstat will print out all the predicted values if you ticked **Fitted Values** when you did the Linear regression

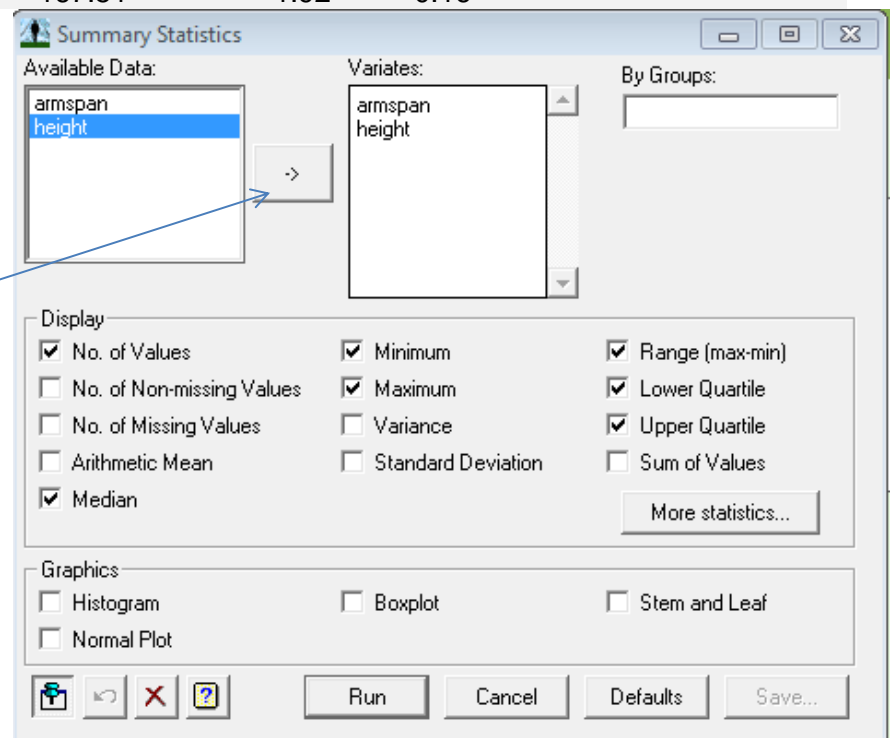
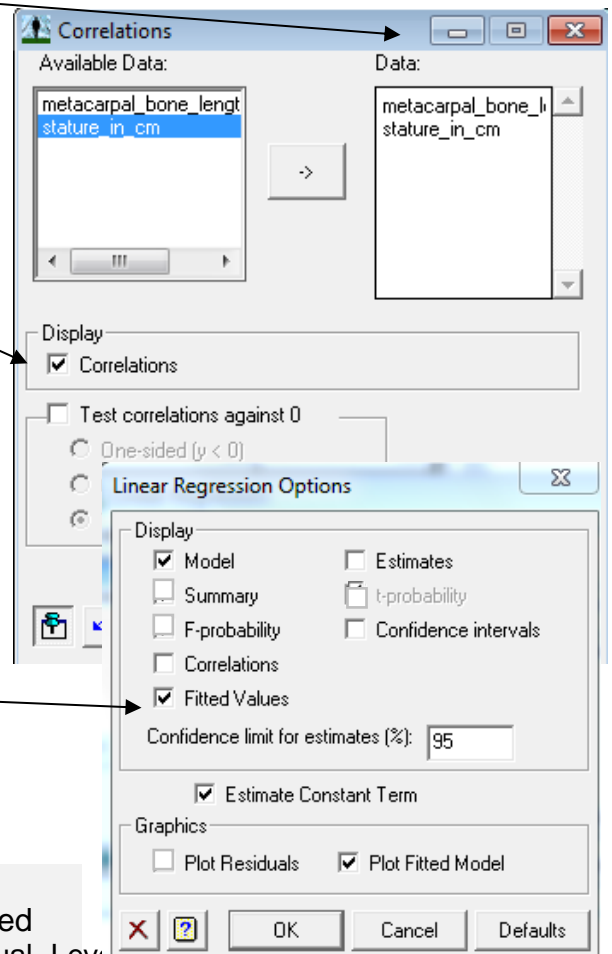
Genstat would have also printed out the standardized residuals if you ticked **Fitted Values**

Fitted values and residuals

Unit	Response	Fitted value	Standardized residual	Level
1	171.00	170.91	0.02	0.11
2	178.00	181.11	-0.92	0.37
3	157.00	160.71	-1.18	0.46
4	163.00	164.11	-0.31	0.28
5	172.00	176.01	-1.03	0.17
6	183.00	177.71	1.40	0.22
7	173.00	172.61	0.10	0.11
8	175.00	167.51	1.92	0.16
9	173.00			
Mean	171.67			

Finding the typical data-using Summary Stats

1. To get Summary statistics, Click on **Summary statistics** from the **Stats** menu
2. Click on the arrow to select what variables you want statistics for



3. Tick the statistics you want to get.
4. Click on Run
5. To see the information, Select **Output** from the **Window** menu

Summary statistics for armspan Summary statistics for height

Number of values = 30
 Median = 157
 Minimum = 134
 Maximum = 167
 Range = 33
 Lower quartile = 150.5
 Upper quartile = 161

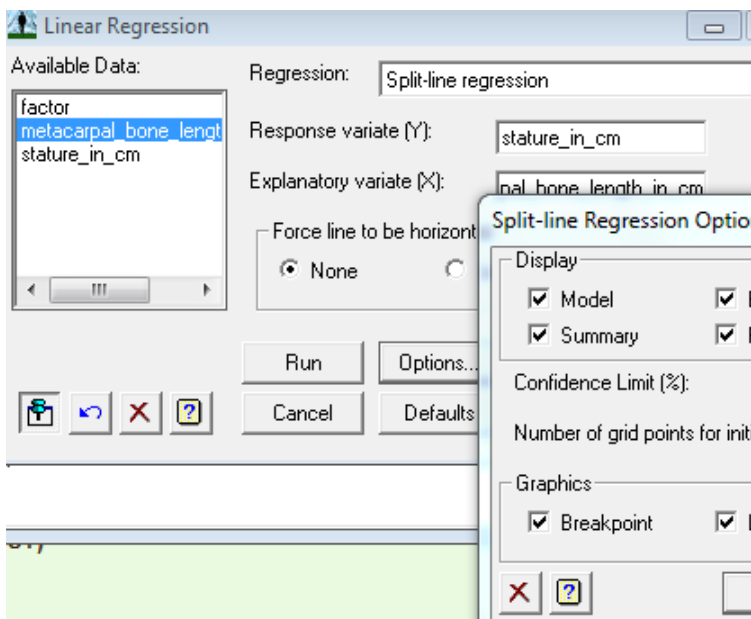
Number of values = 30
 Median = 158
 Minimum = 135
 Maximum = 164
 Range = 29
 Lower quartile = 151.5
 Upper quartile = 161

You can copy these by highlighting clicking on  (or use CTRL and C or right click and use **Copy**).

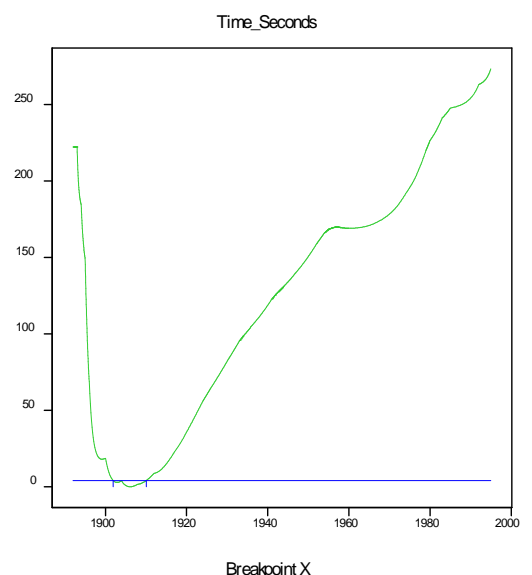
Piecewise Functions

If you think your model would be better as two straight lines rather than one (or even three lines!) you can fit a piecewise model. Genstat will fit the model and even find the best breakpoint (where to split the model) for you.

1. Choose **Stats** menu then **Linear Regression** then change the regression type to **Splitline regression**




As this file is probably best not as piecewise function, you may wish to try with the file And the variables *Mens 1500 m* and use time as the response and year as the explanatory. You will get graphs as shown. You can see that there is a split in the data around 1910. Looking at the output you can see that it is at 1906

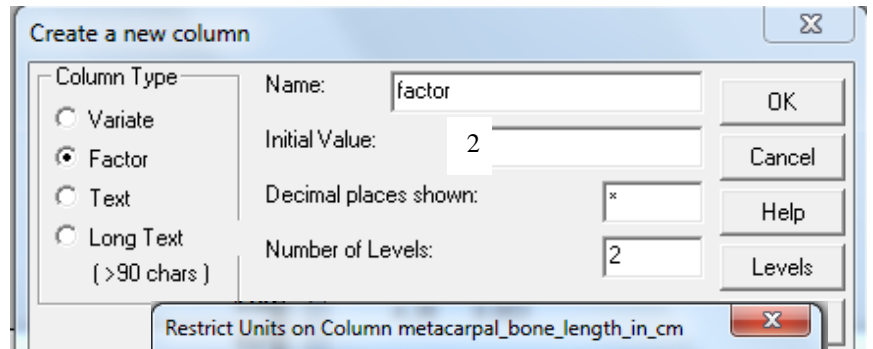


Estimates of parameters

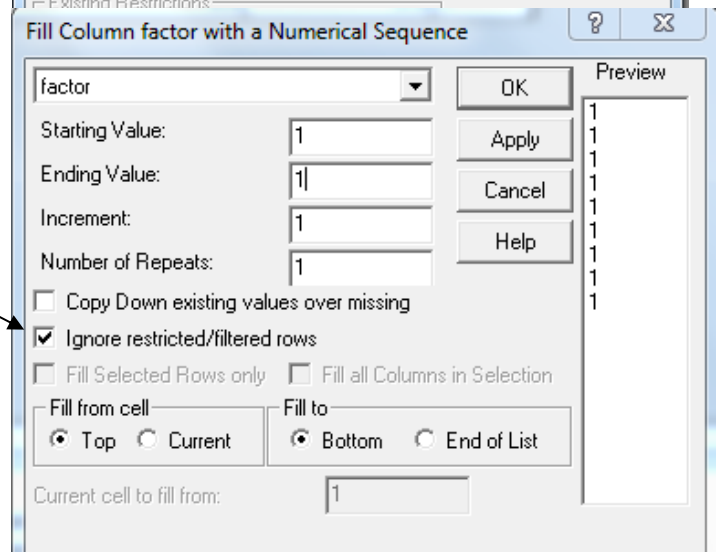
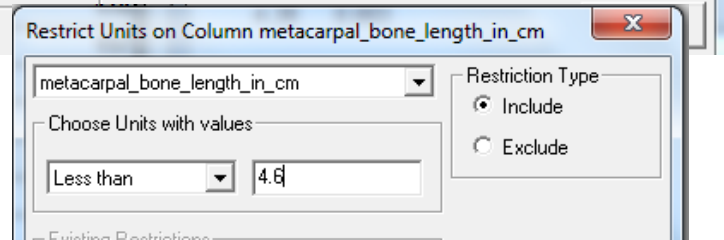
Parameter	estimate	s.e.
Breakpoint_X	1906.07	1.26


However, assume we wish to split the data in the *metacarpal* file at 4.6cm. To do this and graph both models and get the equation for both you will need to divide the data into two groups.

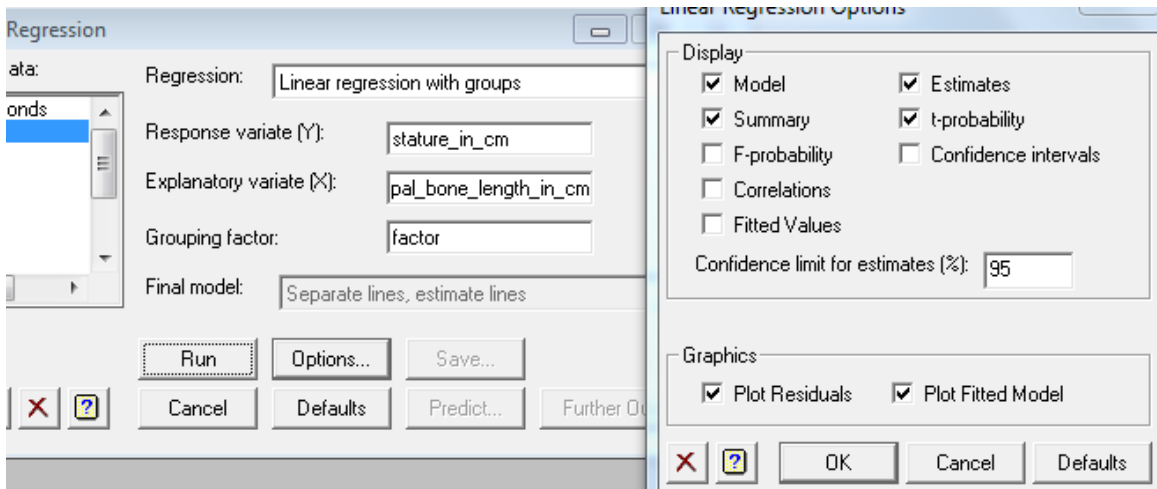
2. Create a factor column  and call it **factor**
3. Go to **Spread** then **Restrict/Filter** then **By value:** - here the data is restricted to all the values where the metacarpal length is less than 4.6



4. From the **Spread** menu, choose **Calculate**, then **Fill** and fill with the value 1 as shown but make sure you tick **Ignore restricted/filtered rows** as shown



5. Remove the filter with 
6. Now you can use **Linear Regression** but use **Linear Regression with groups**

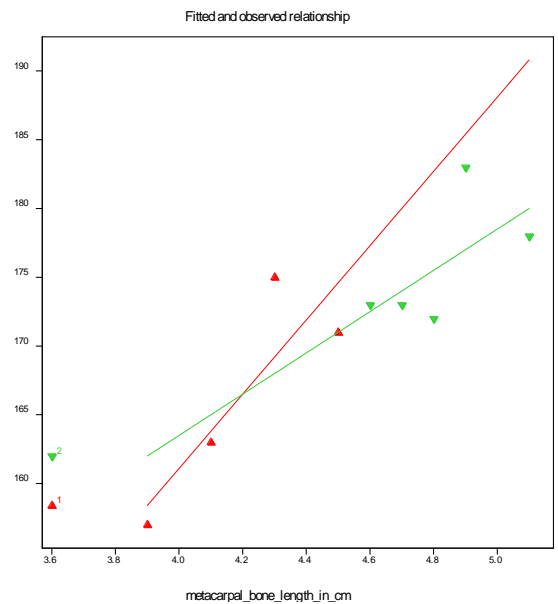


Estimates of parameters

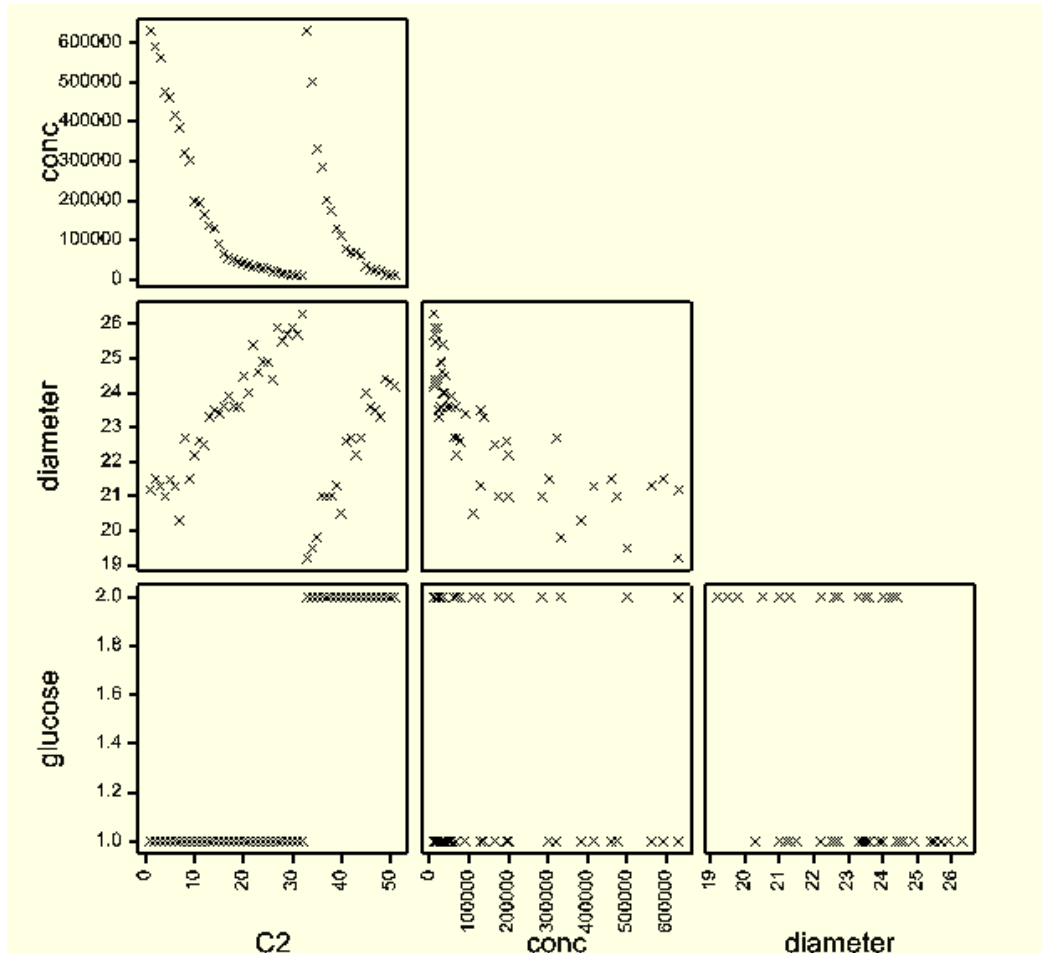
Parameter	estimate	s.e.	t(5)	t pr.
factor 1	53.1	42.6	1.25	0.268
factor 2	103.5	56.8	1.82	0.128
metacarpal_bone_length_in_cm.factor 1	27.0	10.1	2.66	0.045
metacarpal_bone_length_in_cm.factor 2	15.0	11.8	1.27	0.259

While you don't have an r value, you do have the t probabilities and as you can see they are higher than 0.05 and before they were only 0.03 so as mentioned earlier, this data set would be better not as a piecewise model!

Plotting more than one pair of variables at the same time -matrix



When you have more than one pair of data variables, you can plot all the possible data pair combinations by using **Graphics** and then **Scatterplot Matrix** and choosing all the data variables – this now gives you a plot like the one below. Here there are 3 pairs of variables, and the six combinations are plotted – the first row gives X on the x axis and concentration (middle graph) on the y axis and diameter (right graph) on the y axis. The second row has concentration on the x axis with

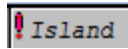


X(left) and diameter (right) on the x axis. The third row has diameter on the x axis and X (left) and concentration (right) on the y axis. This data is from the file cell (but the second column has been deleted – it had 2 values 1 or 2 for glucose)

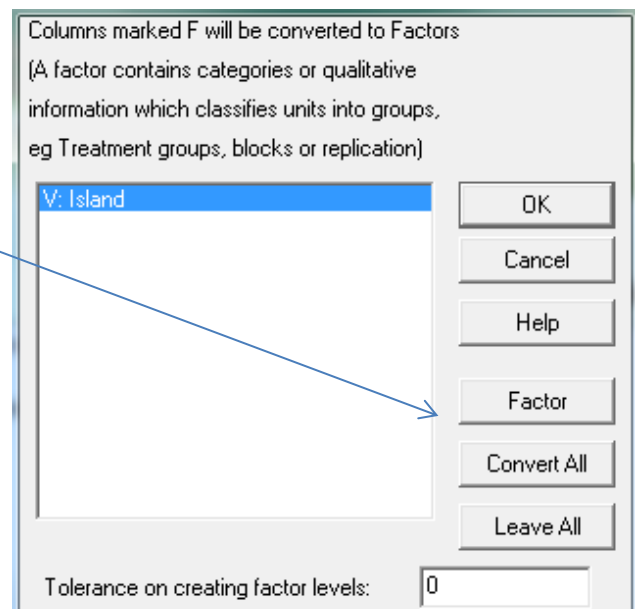
Linear Regression with groups

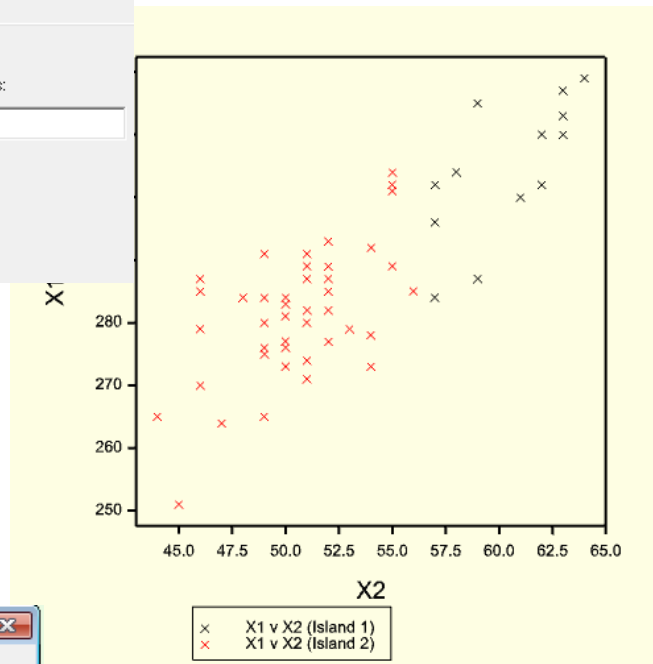
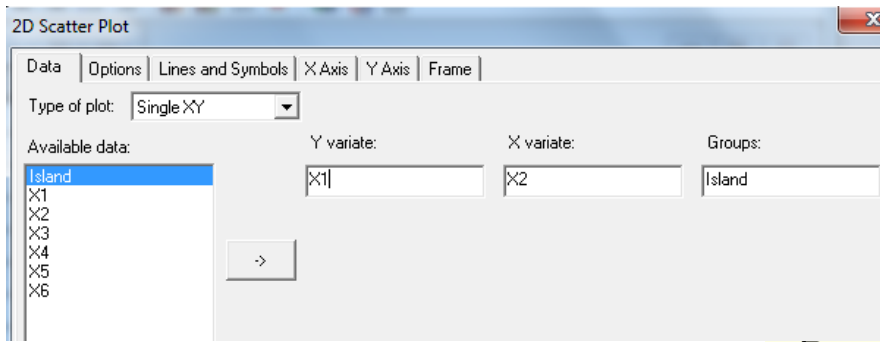
To do this, you must make sure you have your groups as a factor either by

- Clicking on **Factor** or **Convert All** when opening the file
- Or right clicking in the column and choosing **Convert to Factor** when the spreadsheet has been opened

Factor columns have a red exclamation point before the title 

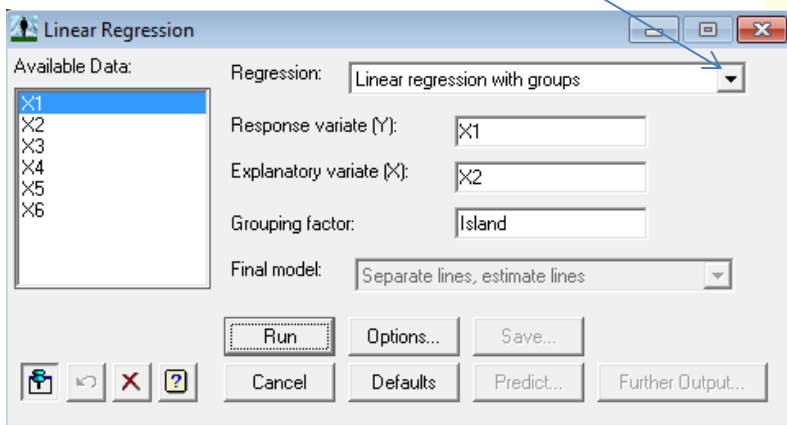
To draw the **2D Scatterplot**, you fill in the dialogue box as shown and each group is shown in a different colour.



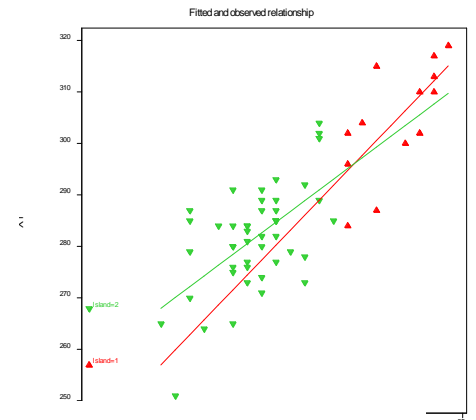


To draw the line of best fit and find the equation for each of the two groups,

1. **Stats, Linear Regression** and then use the arrow to choose **linear regression with groups** and fill in as shown



2. The fitted model shows as two lines of the graph



3. Look under the Here you have 2 gradients and 2 constants, one for each island

So Island 1 equation is:
 $X1 = 2.905 X2 + 129.1$

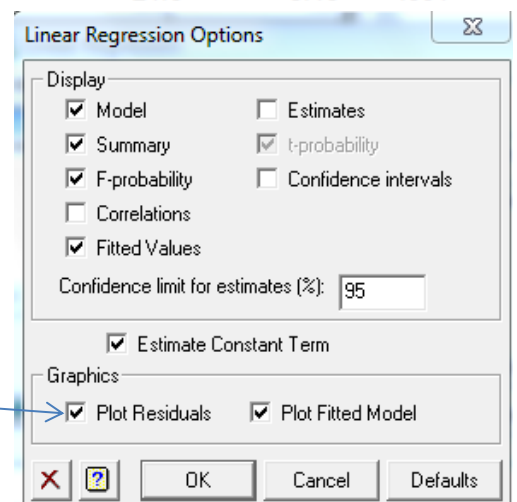
Island 2 equation is:
 $X1 = 2.088 X2 + 176.1$

Estimates of parameters

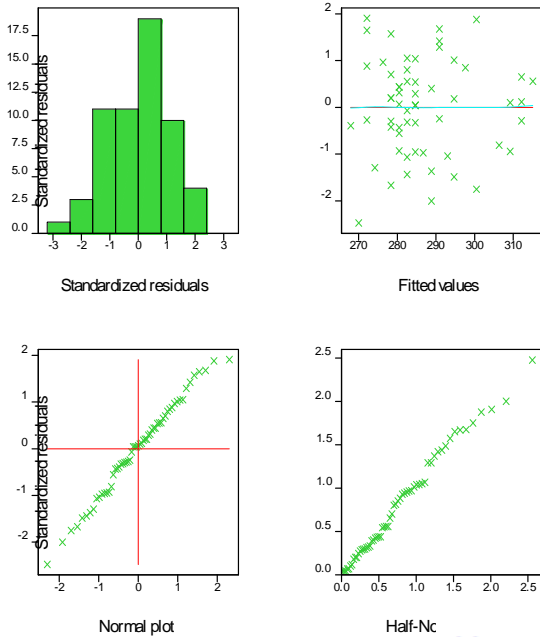
Parameter	estimate	s.e.	t(55)	t pr.
Island 1	129.1	53.9	2.40	0.020
Island 2	176.1	21.6	8.16	<.001
X2.Island 1	2.905			
X2.Island 2	2.088			

Residuals

If you also tick Residual Plot in **Options** when doing the **linear Regression**



X1



Message: the following units have large standardized residuals.

Unit	Response	Residual
27	251.00	-2.48


Outliers

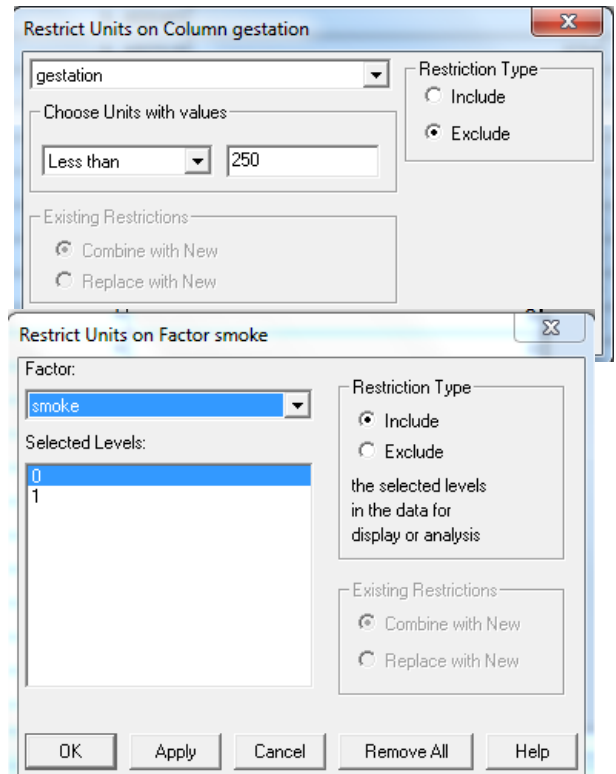
These are identified as shown

- The unit refers to row number in genstat
- The response is the y axis value
- The standardised residual is a measure of how far it is from the line of best fit; anything over 2 or under -2 is considered an outlier. If it negative, it's below the line, if its positive, it's above the line

Restricting/Filtering Data

You can filter by value or by group and then graph or complete Linear Regression on the remaining data.

- Go to **Spread – Restrict Filter**
 - By value.** Set up your restriction ensure you choose between include and exclude correctly. The one opposite will show all gestation lengths *longer* 250 days
 - By group.** Set up your restriction. The one shown will display the group smoking = 0 (non-smokers)
 - Removing the filters.**  will remove filters.



Non- Linear Models

You can fit polynomial, exponential, power, square


Exponential Function

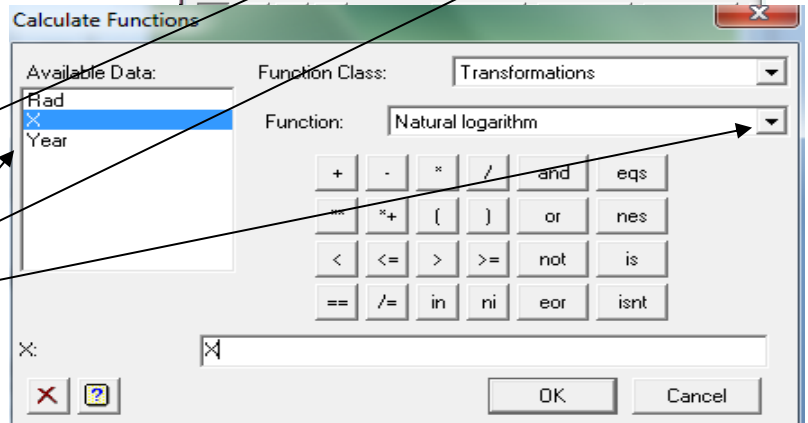
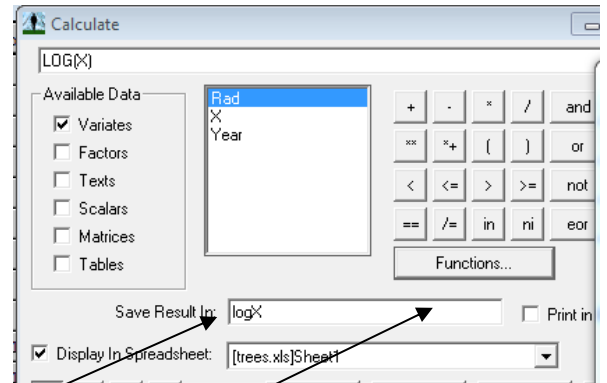
$$y = Ae^{kx} \quad (\text{also can be written } y = Ak^x)$$

e.g. $y = 2e^{3x}$ or $y = 3^x$

Where A is the original amount, r = rate or growth factor, x is time

The file trees has the cross section of a tree trunk. In 1990, when the recording of the cross sections began, the tree trunk which had a cross section of 2cm. Before you can use linear regression you need to transform the data so a linear relationship is present. You can use Natural logarithms to do this.

8. Open the file *trees*. Note: X is the number of years since recording began i.e. 1990.
9. Use the calculator  as before
10. This time we are going to save the results in the spreadsheet.
 - a. Enter in a name for the column of the spreadsheet
 - b. Click on **FUNCTIONS**
 - c. Use the arrow to select Natural logarithm
 - d. Double Click on X
 - e. Click Ok Twice



You will have got a warning message and you can see the new column is highlighted and an * put in Row 1.

Checking the output, there is a warning message

Warning 2, code CA 7, statement 1 on line 66

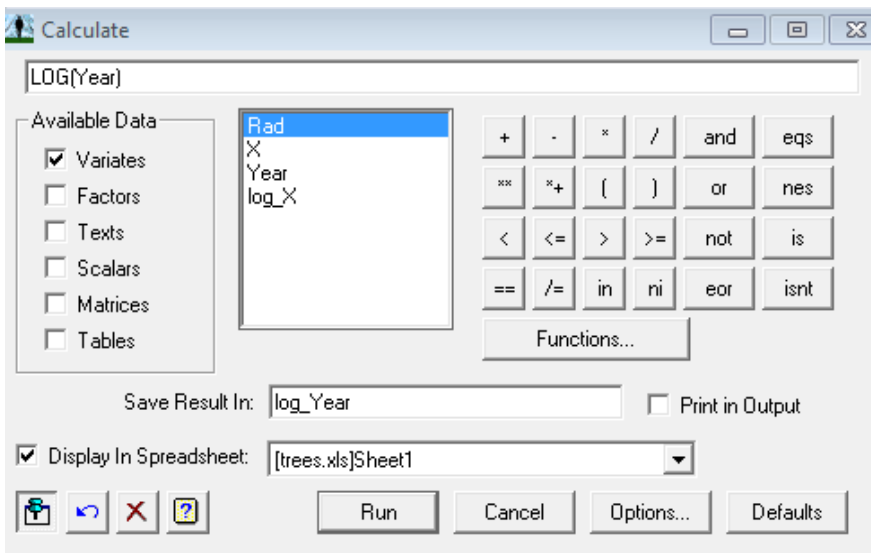
Command: CALCULATE log_X=LOG(X)
Invalid value for argument.

The first argument of the LOG function in unit 1 has the value 0.0000

As you would expect!

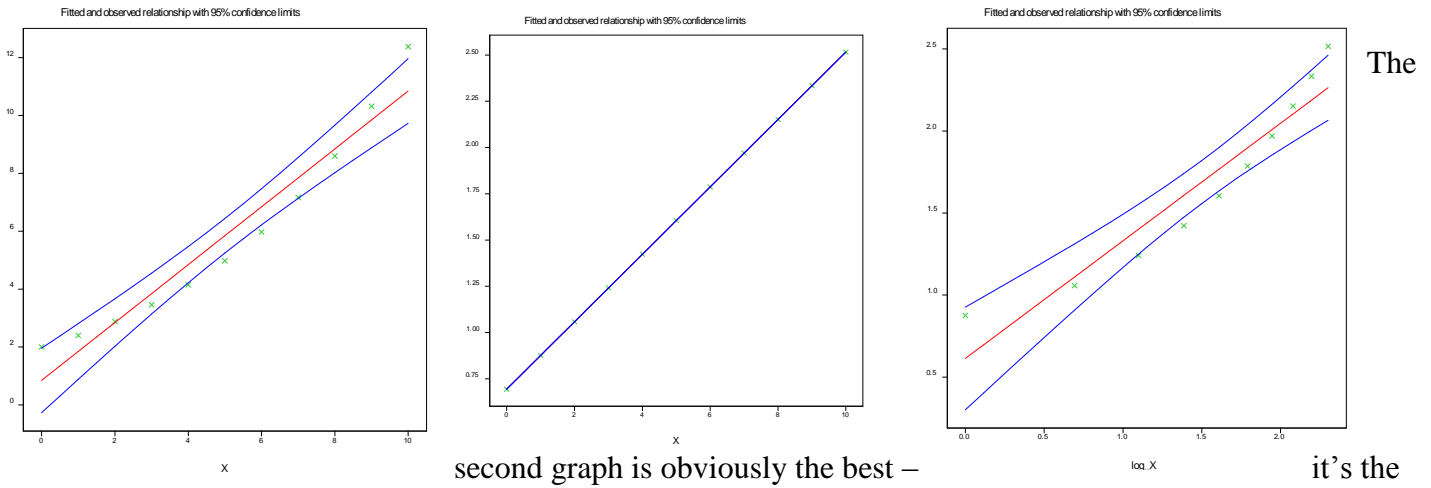
Repeat the transformation for the radius, ensuring you have a new name for the column where the results are to be displayed.

Row	Year	X	Rad	log_X
1	1990	0	2	
2	1991	1	2.4	0
3	1992	2	2.88	0.693147
4	1993	3	3.46	1.09861
5	1994	4	4.15	1.38629
6	1995	5	4.98	1.60944
7	1996	6	5.97	1.79176
8	1997	7	7.17	1.94591
9	1998	8	8.6	2.07944
10	1999	9	10.32	2.19722
11	2000	10	12.38	2.30259



Row	Year	X	Rad	log_X	log_Radius
1	1990	0	2	*	0.693147
2	1991	1	2.4	0	0.875469
3	1992	2	2.88	0.693147	1.05779
4	1993	3	3.46	1.09861	1.24127
5	1994	4	4.15	1.38629	1.42311
6	1995	5	4.98	1.60944	1.60543
7	1996	6	5.97	1.79176	1.78675
8	1997	7	7.17	1.94591	1.96991
9	1998	8	8.6	2.07944	2.15176
10	1999	9	10.32	2.19722	2.33408
11	2000	10	12.38	2.30259	2.51608

- Now perform **Linear Regression** as you have done previously but try different combinations
- c. X explanatory, Radius Response
 - d. X explanatory, log (Radius) Response
 - e. log (X) explanatory, log (Radius) Response

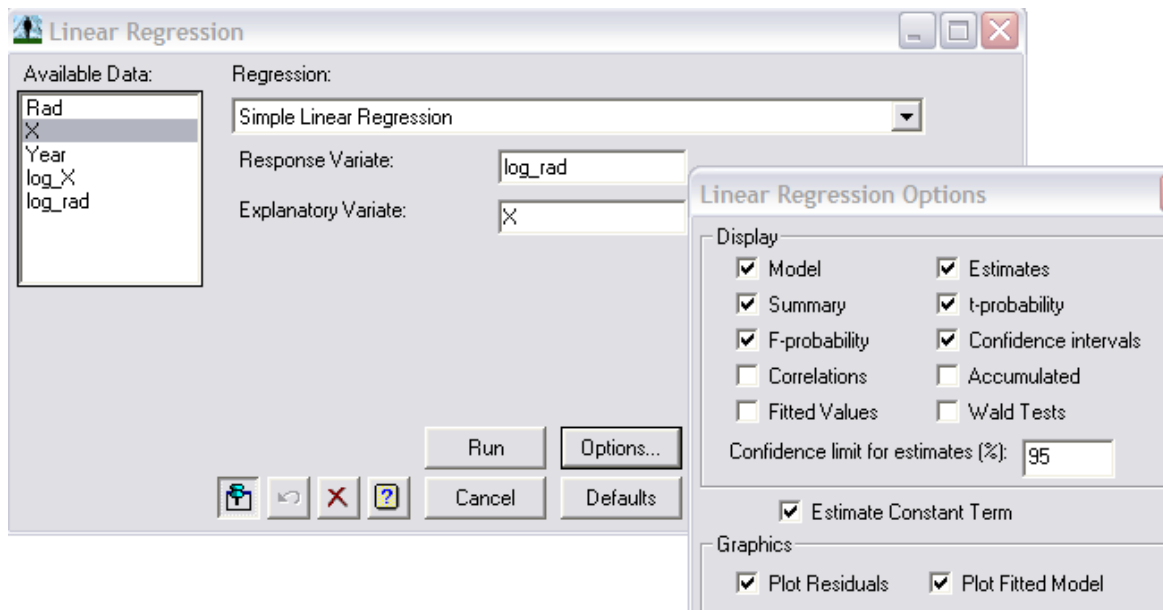


second graph is obviously the best – it's the straightest, also notice the * in the log X column, that's because you cannot log 0, so you cannot use log X to create a model

Year	X	Rad	log_X	log_rad
1990	0	2	*	0.693147
1991	1	2.4	0	0.875469
1992	2	2.88	0.693147	1.05779

This means that an **exponential model** is possibly a very suitable model.

Now you can perform Linear Regression using X as the explanatory variable and log Radius as Response variable as you can see there is a linear relation between the two.



Regression analysis
 Response variate: log_rad
 Fitted terms: Constant, X

Percentage variance accounted for 100.0
 Standard error of observations is estimated to be 0.000475.

Estimates of parameters

Parameter	estimate	s.e.	t(9)	t pr.	lower 95%	upper 95%
Constant	0.693528	0.000268	2589.56	<.001	0.6929	0.6941
X	0.1822906	0.0000453	4026.80	<.001	0.1822	0.1824

Therefore the linear relationship is : **Ln(radius) = 0.1823 x X+ 0.6935**

Transforming this

$$\begin{aligned}
 e^{\text{Ln}(\text{radius})} &= e^{0.1823 \times X + 0.6935} \\
 &= e^{0.1823 \times X} \times e^{0.6935} \\
 \text{radius} &= e^{0.6935} e^{0.1823 \times X} \\
 &= 2.007 e^{0.1823X}
 \end{aligned}$$

We can predict that after seven years, the radius of the tree will be

$$\begin{aligned}
 \text{Radius} &= 2.007 e^{0.1823X} \\
 &= 2.007 e^{0.1823 \times 7} \\
 &= 7.168 \text{ (4sf)}
 \end{aligned}$$

This compares well with the observed value of 7.17.

Power function

$$y=kx^a \quad (\text{e.g. } y=3x^2)$$

A

type of
needs a


Hardener g	5	10	15	20	25	30	35	40
Time taken min	8.8	3.1	1.7	1.1	0.8	0.6	0.5	0.4

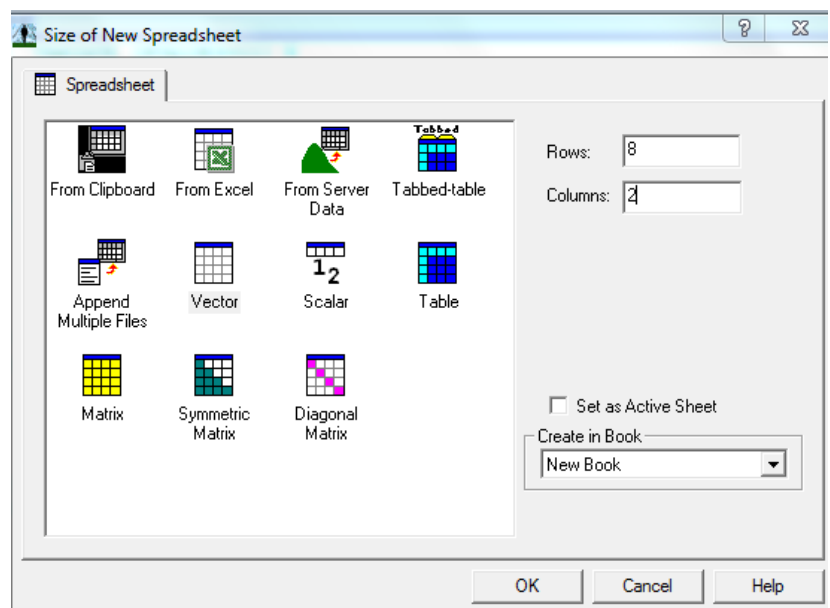
certain
glue
hardener

added to set. The amount of hardener added affects the time taken for the glue to set, as shown in the table above

While this file is available as *glue*, this time we will enter the data in manually.

You may wish to clear the data from the last file first (**Data, Clear All data**)

- g. Click on , you will need 8 rows and 2 columns
- h. Type in the hardener values in the first column and the time taken values in the second column
- i. Right click in the first column and choose **Column Attribute**.



- j. Fill in the dialogue box as shown below. This is where you can also change the type of data by using **Convert** if it is the wrong type (variate when it should be date etc.) and where you can change the **Date Type**. You can alter the width here or by manually dragging in the spreadsheet window.

Column Attributes/Format for C11

Column: C11 Type: [] OK

Name: Hardener_g Variate Cancel

Description: [] Apply

Decimals: * [] Width: 6 Help

Restrict data entered to be in the range: Sheet...

Minimum: * [] Maximum: * [] Convert...

Identifying information used in output: Default Fill..

Date Type...

Justification: Default Left Right Centred

Numeric Format: General Scientific Fixed Date

Column created: 20-Nov-2010 4:43:39 pm

- Repeat for the other column, naming it **Time_taken - min**

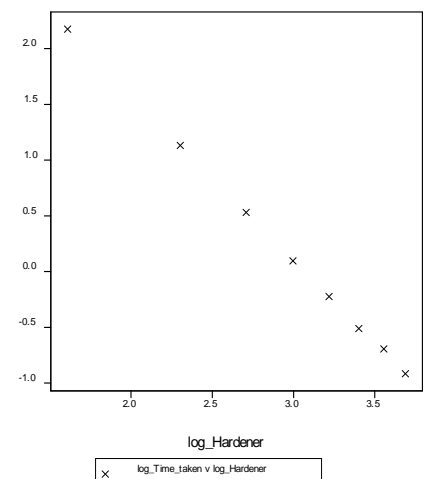
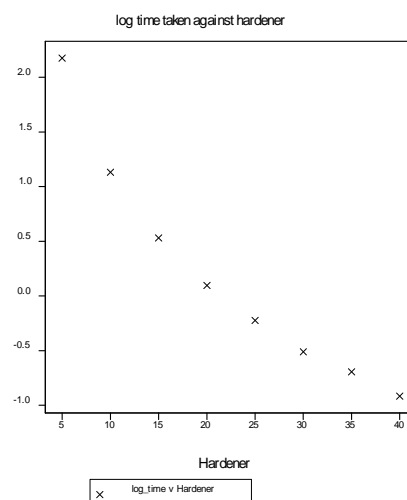
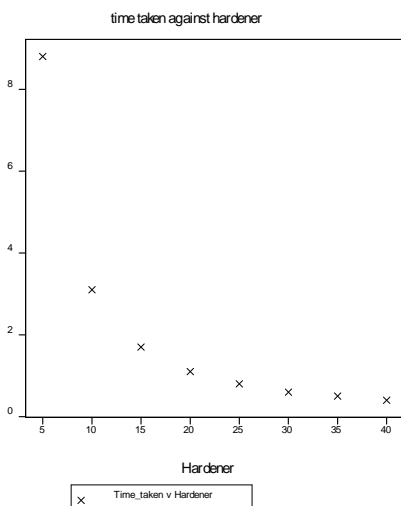
Now you can transform the data as before. (Remember to use **Natural Logarithms**) and graph the three possible models

- Explanatory : Hardener, Response: Time taken
- Explanatory : Hardener, Response: log (Time taken)
- Explanatory : log(Hardener), Response: log(Time taken)

Row	Hardener_g	Time_take	log_Hardener	log_Time_taken
1	5	8.8	1.60944	2.17475
2	10	3.1	2.30259	1.1314
3	15	1.7	2.70805	0.530628
4	20	1.1	2.99573	0.0953102
5	25	0.8	3.21888	-0.223144
6	30	0.6	3.4012	-0.510826

Now graph the three possible models.

The last graph looks the most linear, so perform Linear Regression on Explanatory : log(Hardener), Response: log(Time taken) to find the equation for the power model



regression analysis

Response variate: log_time

Fitted terms: Constant, log_hardener

Estimates of parameters

Parameter	estimate	s.e.	t(6)	t pr.
Constant	4.5504	0.0252	180.42	<.001
log_hardener	-1.48273	0.00838	-176.83	<.001

Parameter	lower95%	upper95%
Constant	4.489	4.612
log_hardener	-1.503	-1.462

In this glue example, the y intercept is 4.55 and the gradient -1.48

$$\ln(\text{hardener}) = -1.48\ln(\text{time}) + 4.55$$

$$e^{\ln(\text{hardener})} = e^{-1.48\ln(\text{time}) + 4.55}$$

$$\text{hardener} = e^{-1.48\ln(\text{time})} \times e^{4.55}$$

$$\text{hardener} = e^{4.55} \times e^{-1.48\ln(\text{time})}$$

$$= 94.6 \text{ time}^{-1.48}$$

$$(-1.48\ln(\text{time}) = \ln(\text{time})^{-1.48})$$

We can test this model to by substituting in a hardener value e.g. 35 and checking the time taken.

$$\text{Time} = 94.6 (35)^{-1.48}$$

$$= 0.49 \text{ very close to the observed } 0.5$$

Now we can use this to predict the time taken for 50g

$$\text{Time} = 99.48 (50)^{-1.5}$$

$$= 0.28 \text{ minutes}$$

Polynomial

You may fit any **polynomial** in Genstat

- Choose Linear Regression but this time change the Regression to Polynomial Regression, then choose whether you want a quadratic, cubic etc, you will get a similar output to before

- **Regression analysis**

-

- Response variate: stature_in_cm

- Fitted terms: Constant + metacarpal_bone_l_length_in_cm

- Submodels: POL(metacarpal_bone_l_length_in_cm; 2)

-

- Standard error of observations is estimated to be 4.16.

-

- *Message: the following units have high leverage.*

Unit	Response	Leverage
2	178.00	0.73
3	157.00	0.74