

Chapter 9 – Re-expressing Data: Get It Straight!

1. Residuals.

- a) The residuals plot shows no pattern. No re-expression is needed.
- b) The residuals plot shows a curved pattern. Re-express to straighten the relationship.
- c) The residuals plot shows a fan shape. Re-express to equalize spread.

2. Residuals.

- a) The residuals plot shows a curved pattern. Re-express to straighten the relationship.
- b) The residuals plot shows a fan shape. Re-express to equalize spread.
- c) The residuals plot shows no pattern. No re-expression is needed.

3. Airline passengers revisited.

- a) The residuals cycle up and down because there is a yearly pattern in the number of passengers departing Oakland, California.
- b) A re-expression should not be tried. A cyclic pattern such as this one cannot be helped by re-expression.

4. Hopkins winds, revisited.

- a) The plot shows a wavy pattern, indicating a pattern that continues year to year as part of an annual cycle.
- b) A re-expression should not be tried. A cyclic pattern such as this one cannot be helped by re-expression.

5. Models.

a) $\ln \hat{y} = 1.2 + 0.8x$

$$\ln \hat{y} = 1.2 + 0.8(2)$$

$$\ln \hat{y} = 2.8$$

$$\hat{y} = e^{2.8} = 16.44$$

b) $\sqrt{\hat{y}} = 1.2 + 0.8x$

$$\sqrt{\hat{y}} = 1.2 + 0.8(2)$$

$$\sqrt{\hat{y}} = 2.8$$

$$\hat{y} = 2.8^2 = 7.84$$

c) $\frac{1}{\hat{y}} = 1.2 + 0.8x$

$$\frac{1}{\hat{y}} = 1.2 + 0.8(2)$$

$$\frac{1}{\hat{y}} = 2.8$$

$$\hat{y} = \frac{1}{2.8} = 0.36$$

d) $\hat{y} = 1.2 + 0.8 \ln x$
 $\hat{y} = 1.2 + 0.8 \ln(2)$
 $\hat{y} = 1.75$

e) $\log \hat{y} = 1.2 + 0.8 \log x$
 $\log \hat{y} = 1.2 + 0.8 \log(2)$
 $\log \hat{y} = 1.440823997\dots$
 $\hat{y} = 10^{1.4408\dots}$
 $\hat{y} = 27.59$

6. More models.

a) $\hat{y} = 1.2 + 0.8 \log x$
 $\hat{y} = 1.2 + 0.8 \log(2)$
 $\hat{y} = 1.44$

b) $\log \hat{y} = 1.2 + 0.8x$
 $\log \hat{y} = 1.2 + 0.8(2)$
 $\log \hat{y} = 2.8$
 $\hat{y} = 10^{2.8} = 630.96$

c) $\ln \hat{y} = 1.2 + 0.8 \ln x$
 $\ln \hat{y} = 1.2 + 0.8 \ln(2)$
 $\ln \hat{y} = 1.7545\dots$
 $\hat{y} = e^{1.7545\dots} = 5.78$

d) $\hat{y}^2 = 1.2 + 0.8x$
 $\hat{y}^2 = 1.2 + 0.8(2)$
 $\hat{y}^2 = 2.8$
 $\hat{y} = \sqrt{2.8} = 1.67$

e) $\frac{1}{\sqrt{\hat{y}}} = 1.2 + 0.8x$
 $\frac{1}{\sqrt{\hat{y}}} = 1.2 + 0.8(2)$
 $\frac{1}{\sqrt{\hat{y}}} = 2.8$
 $\hat{y} = \frac{1}{2.8^2} = 0.128$

7. Gas mileage.

- a) The association between weight and fuel efficiency of cars is fairly linear, strong, and negative. Heavier cars tend to have lower fuel efficiency.
- b) For each additional thousand pounds of weight, the linear model predicts a decrease of 7.652 miles per gallon in gas mileage.
- c) The linear model is not appropriate. There is a curved pattern in the residuals plot. The model tends to underestimate gas mileage for cars with relatively low and high gas mileages, and overestimates the gas mileage of cars with average gas mileage.

8. Crowdedness.

- a) The scatterplot shows that the relationship between Crowdedness and GDP is strong, negative, and curved. Re-expression may yield an association that is more linear.

- b) Start with logs, since GDP is non-negative. A plot of the log of GDP against Crowdedness score may be straighter.

9. Gas mileage revisited.

- a) The residuals plot for the re-expressed relationship is much more scattered. This is an indication of an appropriate model.
- b) The linear model that predicts the number of gallons per 100 miles in fuel efficiency from the weight of a car is: $\widehat{Gal/100} = 0.625 + 1.178(Weight)$.
- c) For each additional 1000 pounds of weight, the model predicts that the car will require an additional 1.178 gallons to drive 100 miles.
- d)

$$\widehat{Gal/100} = 0.625 + 1.178(Weight)$$

$$\widehat{Gal/100} = 0.625 + 1.178(3.5)$$

$$\widehat{Gal/100} = 4.748$$

According to the model, a car that weighs 3500 pounds (3.5 thousand pounds) is expected to require approximately 4.748 gallons to drive 100 miles, or 0.04748 gallons per mile.

This is $\frac{1}{0.04748} \approx 21.06$ miles per gallon.

10. Crowdedness again.

- a) This re-expression is not useful. The student has gone too far down the ladder of powers. We now see marked downward curvature and increasing scatter.
- b) The next step would be to try a “weaker” re-expression, like reciprocal square root or log of GDP. Having gone too far, the student should move back “up” the ladder of powers.

11. GDP.

- a) Although more than 87.6% of the variation in GDP can be accounted for by the model, the residuals plot should be examined to determine whether or not the model is appropriate.
- b) This is not a good model for these data. The residuals plot shows curvature.

12. Treasury bills.

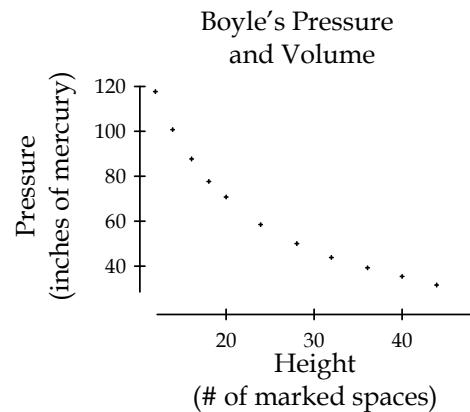
Re-expression should not be tried. An erratic trend that is positive then negative cannot be straightened by re-expression.

13. Better GDP model?

This is not a better model. The residuals plot still has a strong pattern.

14. Pressure.

The scatterplot at the right shows a strong, curved, negative association between the height of the cylinder and the pressure inside. Because of the curved nature of the association, a linear model is not appropriate.



Re-expressing the pressure as the reciprocal of the pressure produces a scatterplot that is much straighter. Computer regression output for the height versus the reciprocal of pressure is below.

Dependent variable is: **recip pressure**

No Selector

R squared = 100.0% R squared (adjusted) = 100.0%

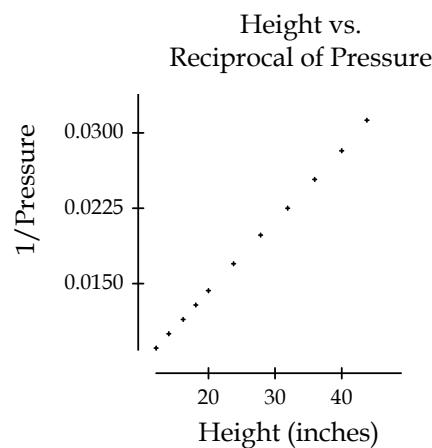
s = 0.0001 with 12 - 2 = 10 degrees of freedom

| Source | Sum of Squares | df | Mean Square | F-ratio |
|------------|----------------|----|-------------|---------|
| Regression | 0.000841 | 1 | 0.000841 | 75241 |
| Residual | 0.000000 | 10 | 0.000000 | |

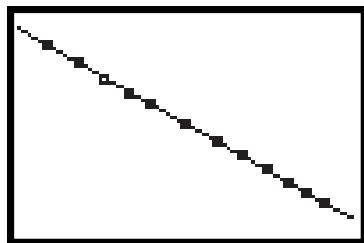
| Variable | Coefficient | s.e. of Coeff | t-ratio | prob |
|----------|-------------|---------------|---------|----------|
| Constant | -7.66970e-5 | 0.0001 | -0.982 | 0.3494 |
| Height | 7.13072e-4 | 0.0000 | 274 | < 0.0001 |

The reciprocal re-expression is very straight (perfectly straight, as far as the statistical software is concerned!). $R^2 = 100\%$, meaning that 100% of the variability in the reciprocal of pressure is explained by the model. The equation of the model is:

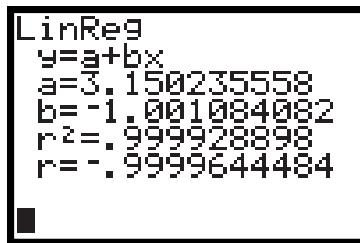
$$\widehat{\frac{1}{\text{Pressure}}} = -0.000077 + 0.000713(\text{Height}).$$



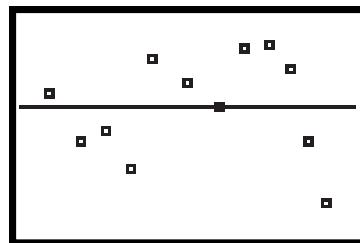
Re-expressing each variable using logarithms is also a good model. TI-83 regression output for the log-log re-expression is given below.



Scatterplot of Log(pressure)
vs. Log(height)



Regression Output



Residuals Plot

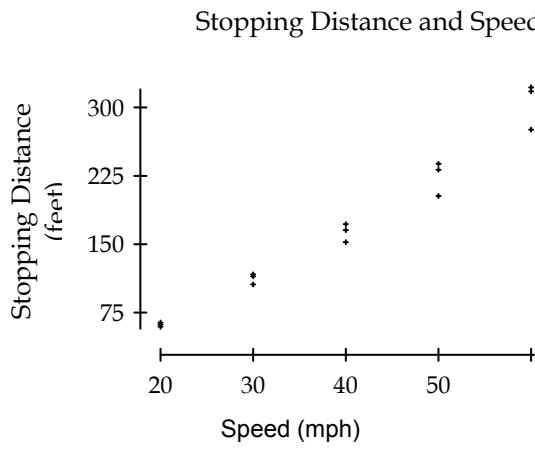
There is a strong, negative, linear relationship between $\log(Pressure)$ and $\log(Height)$.

$\widehat{\log(Pressure)} = 3.150 - 1.001(\log(Height))$ models the situation well, explaining nearly 100% of the variability in the logarithm of pressure. The residuals plot is fairly scattered, indicating an appropriate model.

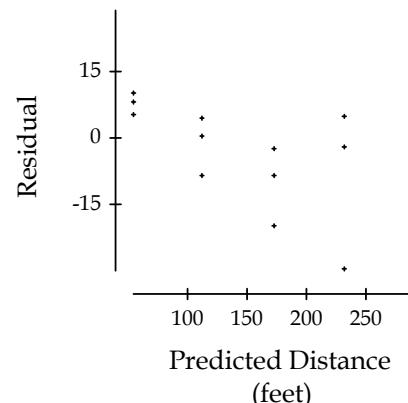
15. Brakes.

- a) The association between speed and stopping distance is strong, positive, and appears straight. Higher speeds are generally associated with greater stopping distances. The linear regression model, with equation

$\widehat{Distance} = -65.9 + 5.98(Speed)$, has $R^2 = 96.9\%$, meaning that the model explains 96.9% of the variability in stopping distance. However, the residuals plot has a curved pattern. The linear model is not appropriate. A model using re-expressed variables should be used.

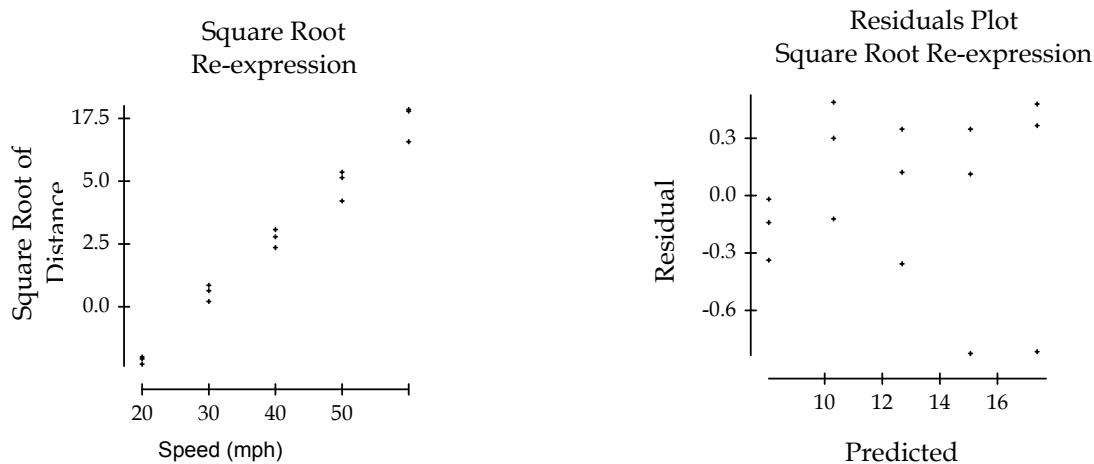


Residuals Plot



158 Part II Exploring Relationships Between Variables

- b) Stopping distances appear to be relatively higher for higher speeds. This increase in the rate of change might be able to be straightened by taking the square root of the response variable, stopping distance. The scatterplot of Speed versus $\sqrt{Distance}$ seems like it might be a bit straighter.



- c) The model for the re-expressed data is $\widehat{\sqrt{Distance}} = 3.303 + 0.235(Speed)$. The residuals plot shows no pattern, and $R^2 = 98.4\%$, so 98.4% of the variability in the square root of the stopping distance can be explained by the model.

d)

$$\widehat{\sqrt{Distance}} = 3.303 + 0.235(Speed)$$

$$\widehat{\sqrt{Distance}} = 3.303 + 0.235(55)$$

$$\widehat{\sqrt{Distance}} = 16.228$$

$$\widehat{Distance} = 16.228^2 \approx 263.4$$

According to the model, a car traveling 55 mph is expected to require approximately 263.4 feet to come to a stop.

e)

$$\widehat{\sqrt{Distance}} = 3.303 + 0.235(Speed)$$

$$\widehat{\sqrt{Distance}} = 3.303 + 0.235(70)$$

$$\widehat{\sqrt{Distance}} = 19.753$$

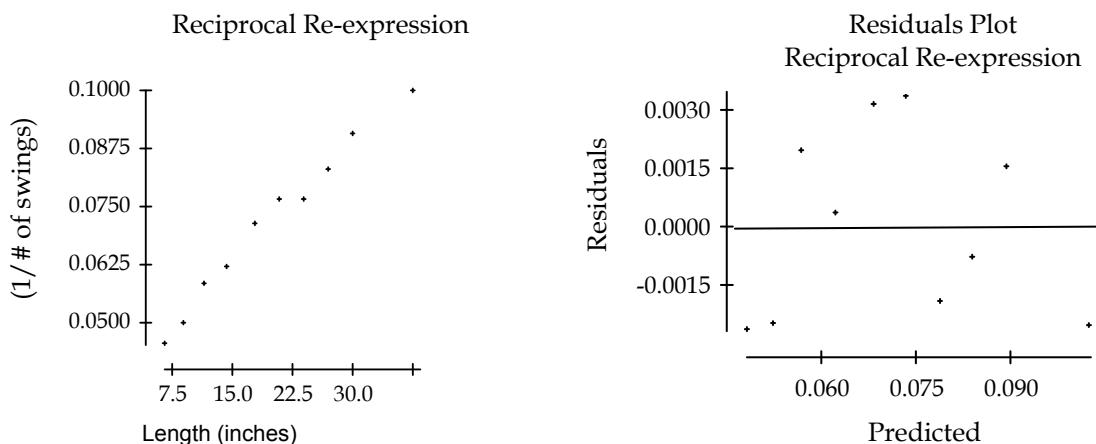
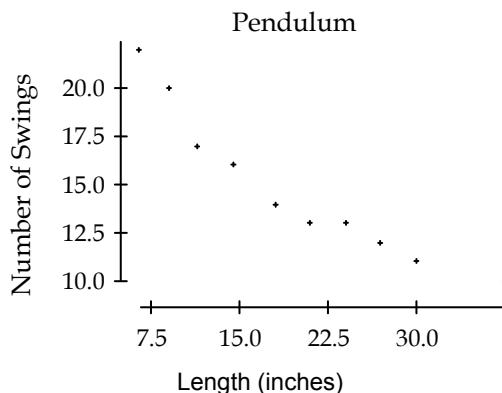
$$\widehat{Distance} = 19.753^2 \approx 390.2$$

According to the model, a car traveling 70 mph is expected to require approximately 390.2 feet to come to a stop.

- f) The level of confidence in the predictions should be quite high. R^2 is high, and the residuals plot is scattered. The prediction for 70 mph is a bit of an extrapolation, but should still be reasonably close.

16. Pendulum.

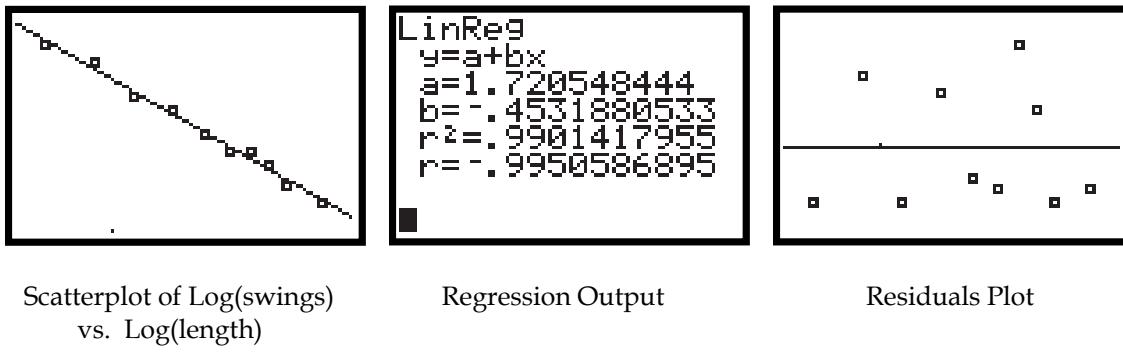
- a) The scatterplot shows the association between the length of string and the number of swings a pendulum took every 20 seconds to be strong, negative, and curved. A pendulum with a longer string tended to take fewer swings in 20 seconds. The linear model is not appropriate, because the association is curved.
- b) Curvature in a negative relationship sometimes is an indication of a reciprocal relationship. Try re-expressing the response variable with the reciprocal.



- c) The reciprocal re-expression yields the model $\widehat{\frac{1}{Swings}} = 0.0367 + 0.00176(Length)$.

The residuals plot is scattered, and $R^2 = 98.1\%$, indicating that the model explains 98.1% of the variability in the reciprocal of the number of swings. The model is both appropriate and accurate.

Re-expressing each variable with logarithms also results in a good model, and has the added benefit of maintaining the direction of the relationship between number of swings and the length of the string. TI-83 regression output for the log-log re-expression is given below.



There is a strong, negative, linear relationship between log(swings) and log(length).

$\widehat{\log(\text{Swings})} = 1.721 - 0.453(\log(\text{Length}))$ explains 99% of the variability in log(swings), and the residuals plot is scattered. The log-log re-expression is useful, as well.

d) Using the reciprocal model:

$$\widehat{\frac{1}{\text{Swings}}} = 0.0367 + 0.00176(\text{Length}) = 0.0367 + 0.00176(4) = 0.04374$$

$$\widehat{\text{Swings}} = \frac{1}{0.04374} \approx 22.9$$

According to the reciprocal model, a pendulum with a 4" string is expected to swing approximately 22.9 times in 20 seconds.

Using the log-log model:

To make estimates involving logs, it is a good idea to use as much accuracy as possible. If the equation of the model is stored in the calculator, the model predicts that a pendulum with a 4" string will swing approximately 28 times in 20 seconds.

e) Using the reciprocal model:

$$\widehat{\frac{1}{\text{Swings}}} = 0.0367 + 0.00176(\text{Length}) = 0.0367 + 0.00176(48) = 0.12118$$

$\widehat{\text{Swings}} = \frac{1}{0.12118} \approx 8.3$. The model predicts 8.3 swings in 20 seconds for a 48" string.

Using the log-log model:

To make estimates involving logs, it is a good idea to use as much accuracy as possible. If the equation of the model is stored in the calculator, the model predicts that a pendulum with a 48" string will swing approximately 9.1 times in 20 seconds.

- f) Confidence in the predictions is fairly high. The model is appropriate, as indicated by the scattered residuals plot, and accurate, indicated by the high value of R^2 . The only concern is the fact that these predictions are slight extrapolations. The lengths of the strings aren't too far outside the range of the data, so the predictions should be reasonably accurate.

17. Baseball salaries 2012.

- a) The scatterplot of year versus highest salary is moderately strong, positive and curved. This is not surprising, since inflation is exponential. If the highest baseball salaries have been rising with inflation, we would expect them to be increasing exponentially as well.
- b) The linear model is not an appropriate model. The residuals plot shows a strong bend. We should try re-expressing the data.
- c) Re-expression using the natural logarithm was successful since the residuals plot is scattered. According to the model, the highest salary in 2015 is predicted to be:

$$\widehat{\ln(\text{Salary})} = -251.28738 + 0.1270045(\text{Year})$$

$$\widehat{\ln(\text{Salary})} = -251.28738 + 0.1270045(2015)$$

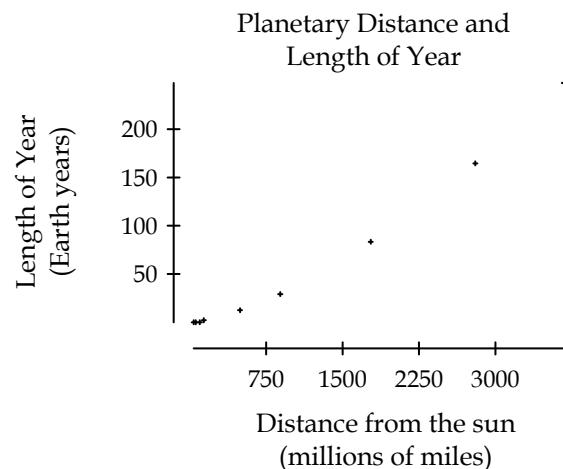
$$\widehat{\ln(\text{Salary})} = 4.6266875$$

$$\widehat{\text{Salary}} = e^{4.6266875} \approx \$102.2 \text{ million}$$

- d) The next three salaries do not fit the previous pattern. The prediction for 2015 was certainly an extrapolation, and these new data points provide evidence that it may not be a reasonable estimate for the highest salary in 2015. None of the three more recent salaries even exceed A-Rod's 2005 salary.

18. Planet distances and years 2012.

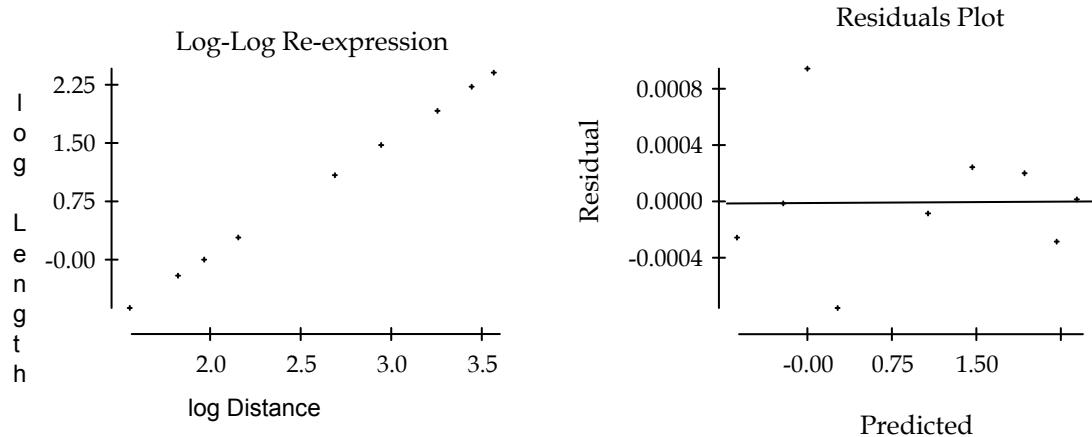
- a) The association between distance from the sun and planet year is strong, positive, and curved concave upward. Generally, planets farther from the sun have longer years than closer planets.



- b) The rate of change in length of year per unit distance appears to be increasing, but not exponentially. Re-expressing with the logarithm of each variable may straighten a plot such as this. The scatterplot and residuals plot for the linear model relating $\log(\text{Distance})$ and $\log(\text{Length of Year})$ appear below.

The regression model for the log-log re-expression is :

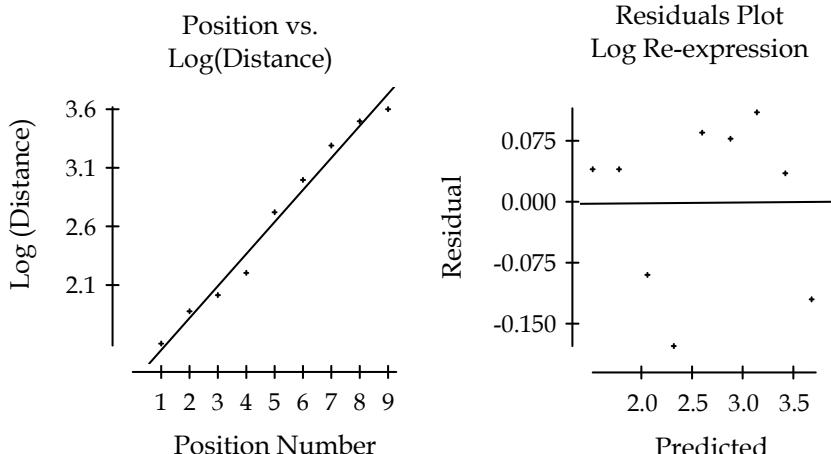
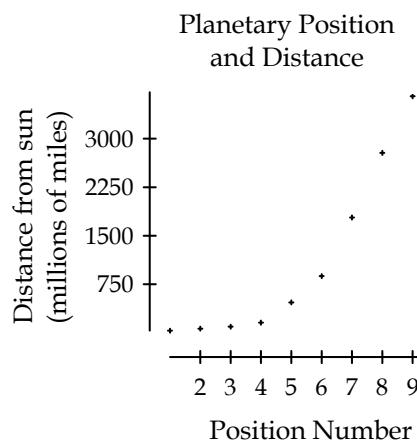
$$\widehat{\log(\text{Length})} = -2.95 + 1.5(\log(\text{Distance})).$$



- c) $R^2 = 100\%$, so the model explains 100% of the variability in the log of the length of the planetary year, at least according to the accuracy of the statistical software. The residuals plot is scattered, and the residuals are all extremely small. This is a very accurate model.

19. Planet distances and order 2012.

- a) The association between planetary position and distance from the sun is strong, positive, and curved (shown at the right). A good re-expression of the data is position versus $\text{Log}(\text{Distance})$. The scatterplot with regression line (below left) shows the straightened association. The equation of the model is $\log(\hat{\text{Distance}}) = 1.245 + 0.271(\text{Position})$. The residuals plot (below right) may have some pattern, but after trying several re-expressions, this is the best that can be done. $R^2 = 98.2\%$, so the model explains 98.2% of the variability in the log of the planets distance from the sun.



- b) At first glance, this model appears to provide little evidence to support the contention of the astronomers. Pluto appears to fit the pattern, although Pluto's distance from the sun is a bit less than expected. A model generated without Pluto does not have a dramatically improved residuals plot, does not have a significantly higher R^2 , nor a different slope. Pluto does not appear to be influential.

But don't forget that a logarithmic scale is being used for the vertical axis. The higher up the vertical axis you go, the greater the effect of a small change.

$$\log(\hat{\text{Distance}}) = 1.24418 + 0.271229(\text{Position})$$

$$\log(\hat{\text{Distance}}) = 1.24418 + 0.271229(9)$$

$$\log(\hat{\text{Distance}}) = 3.685241$$

$$\hat{\text{Distance}} = 10^{3.685241} \approx 4844$$

According to the model, the 9th planet in the solar system is predicted to be approximately 4844 million miles away from the sun. Pluto is actually 3707 million miles away.

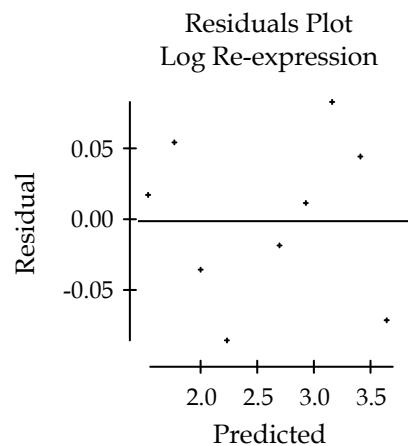
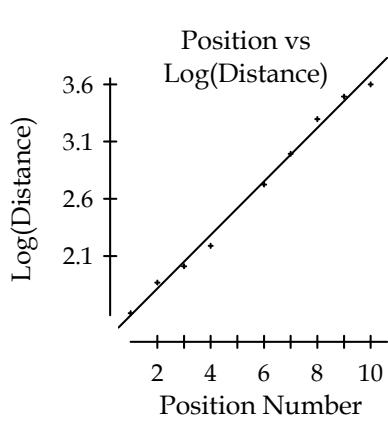
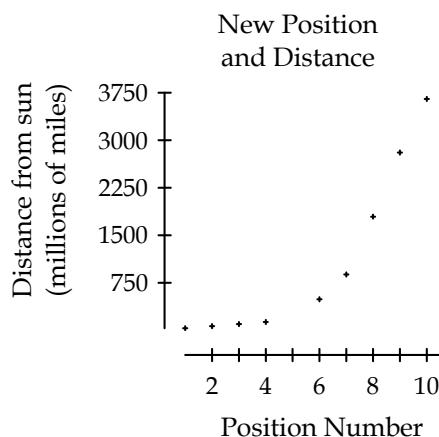
Pluto doesn't fit the pattern for position and distance in the solar system. In fact, the model made with Pluto included isn't a good one, because Pluto influences those predictions. The model without Pluto,

$\log(\overline{\text{Distance}}) = 1.20267 + 0.283680(\text{Position})$, works much better. It has a high R^2 , and scattered residuals plot. This new model predicts that the 9th planet should be a whopping 5699 million miles away from the sun! There is evidence that the IAU is correct. Pluto doesn't behave like planet in its relation to position and distance.

20. Planets 2012, part 3.

Using the revised planetary numbering system, and straightening the scatterplot using the same methods as in Exercise 19, the new model,

$\log(\overline{\text{Distance}}) = 1.32 + 0.23(\text{Position})$, is a slightly better fit. The residuals plot is more scattered, and R^2 is slightly higher, with the improved model explaining 99.5% of the variability in the log of distance from the sun.



Pluto still doesn't fit very well. The new model predicts that Pluto, as 10th planet, should be about 4169 million miles away. That's about 462 million miles farther away. A better model yet is $\log(\overline{\text{Distance}}) = 1.28514 + 0.238826(\text{Position})$, a model made with the new numbering system and with Pluto omitted.

21. Eris: Planets 2012, part 4.

A planet ninth from the sun was predicted, in a previous exercise, to be about 4844 million miles away from the sun. This distance is much shorter than the actual distance of Eris, about 6300 miles.

22. Models and laws: Planets 2012 part 5

The re-expressed data relating distance and year length are better described by their model than the re-expressed data relating position and distance. The model relating distance and year length has $R^2 = 100\%$, and a very scattered residuals plot (with minuscule residuals), possibly a natural “law”. If planets in another solar system followed the Titius-Bode pattern, this belief would be reinforced. Similarly, if data were acquired from planets in another solar system that did not follow this pattern, we would be unlikely to think that this relationship was a universal law.

23. Logs (not logarithms).

- a) The association between the diameter of a log and the number of board feet of lumber is strong, positive, and curved. As the diameter of the log increases, so does the number of board feet of lumber contained in the log.

The model used to generate the table used by the log buyers is based upon a square root re-expression. The values in the table correspond exactly to the model

$$\sqrt{BoardFeet} = -4 + Diameter .$$

b)

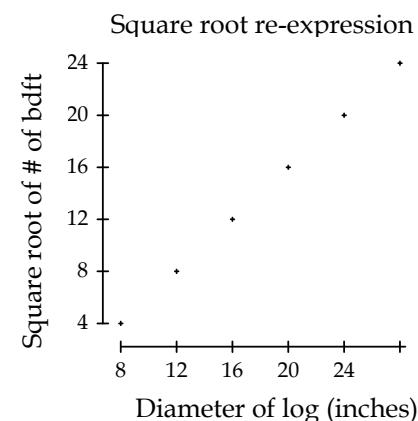
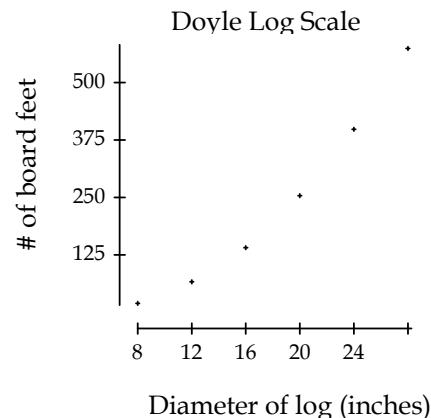
$$\widehat{\sqrt{BoardFeet}} = -4 + Diameter$$

$$\widehat{\sqrt{BoardFeet}} = -4 + (10)$$

$$\widehat{\sqrt{BoardFeet}} = 6$$

$$\widehat{BoardFeet} = 36$$

According to the model, a log 10" in diameter is expected to contain 36 board feet of lumber.



c)

$$\widehat{\sqrt{BoardFeet}} = -4 + Diameter$$

$$\widehat{\sqrt{BoardFeet}} = -4 + (36)$$

$$\widehat{\sqrt{BoardFeet}} = 32$$

$$\widehat{BoardFeet} = 1024$$

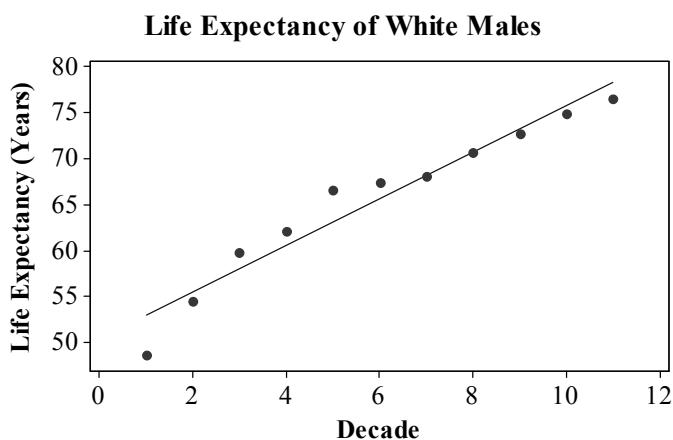
According to the model, a log 36" in diameter is expected to contain 1024 board feet of lumber.

Normally, we would be cautious of this prediction, because it is an extrapolation beyond the given data, but since this is a prediction made from an exact model based on the volume of the log, the prediction will be accurate.

24. Weightlifting 2012

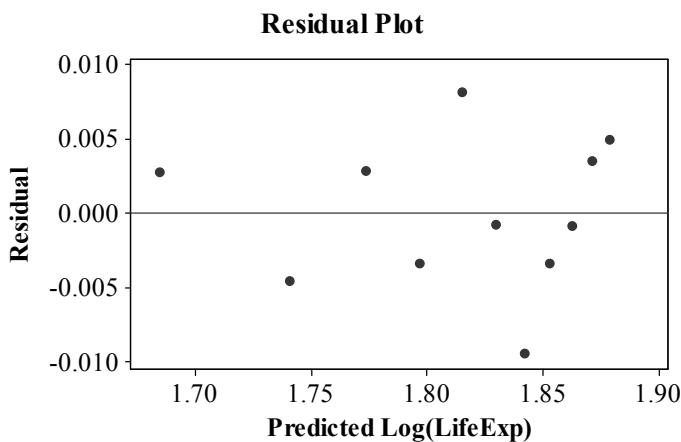
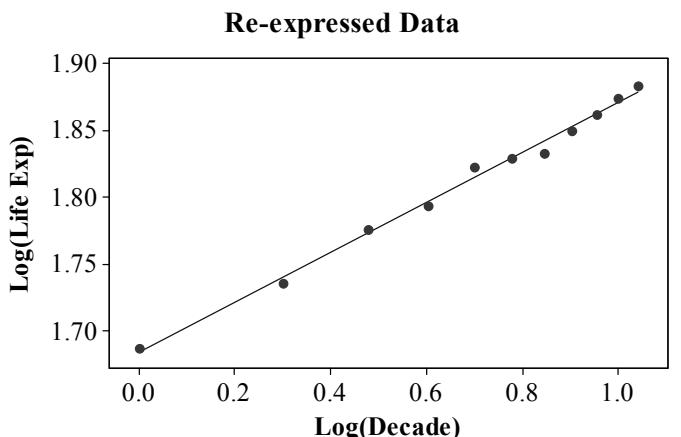
- a) The residual plot for the linear regression between weight class and weight lifted for gold medal winners in weightlifting at the 2012 Olympics shows a curved pattern, indicating that the linear model has failed to model the association well. A re-expressed model might fit the association between weight class and weight lifted better than the linear model.
- b) Both residuals plots are still curved, indicating that neither model is a good one. In the absence of a good model, the model built from re-expressing both weight class and weight lifted with logarithms is the preferable one, since it has a higher value of R^2 .

25. Life expectancy.



The association between year and life expectancy is strong, curved and positive. As the years passed, life expectancy for white males has increased. Re-expressing both variables with logarithms straightens the association significantly.

The re-expressed model, $\widehat{\text{Log(LifeExp)}} = 1.684 + 0.187077 \text{Log(Decade)}$ has a scattered residuals plot and high R^2 , so it will be a good model to predict future increases in life expectancy, as long as we don't attempt to predict too far into the future. The pattern may not continue.



26. Lifting record weight 2012.

a) The reciprocal model is: $\widehat{\frac{1}{WtLifted}} = 0.00427 - 0.000019006(WtClass)$.

$$\widehat{\frac{1}{WtLifted}} = 0.00427 - 0.000019006(77)$$

$$\widehat{\frac{1}{WtLifted}} = 0.002806538$$

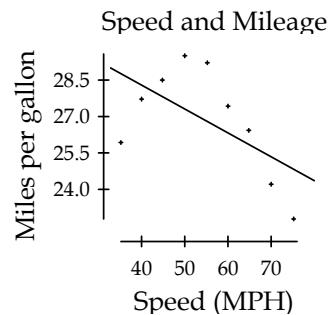
$$\widehat{WtLifted} = 356.31$$

According to the reciprocal model, the winner of the 77 kg weight class was predicted to lift 356.31 kg. Since Xiaojun actually lifted 379 kg, his residual was $379 - 356.31 = 22.69$ kg. He lifted 22.69 kg more than predicted.

b) It is not surprising that a world record lift had a positive residual.

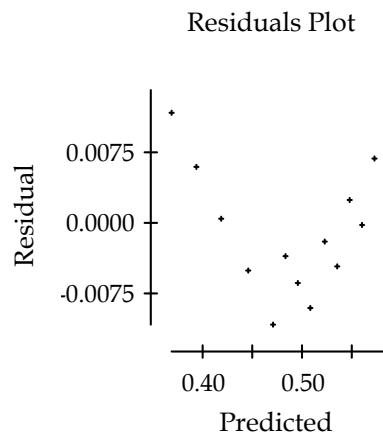
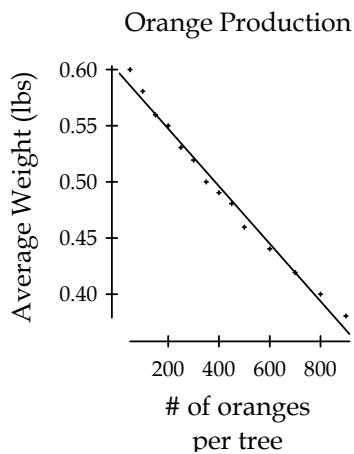
27. Slower is cheaper?

The scatterplot shows the relationship between speed and mileage of the compact car. The association is extremely strong and curved, with mileage generally increasing as speed increases, until around 50 miles per hour, then mileage tends to decrease as speed increases. The linear model is a very poor fit, but the change in direction means that re-expression cannot be used to straighten this association.

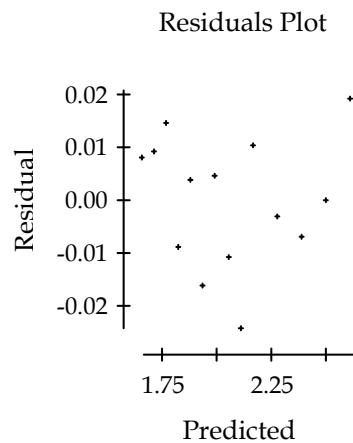
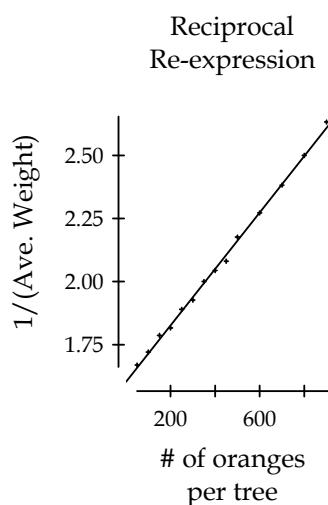


28. Orange production.

The association between the number of oranges per tree and the average weight is strong, negative, and appears linear at first look. Generally, trees that contain larger numbers of oranges have lower average weight per orange. The residuals plot shows a strong curved pattern. The data should be re-expressed.



Plotting the number of oranges per tree and the reciprocal of the average weight per orange straightens the relationship considerably. The residuals plot shows little pattern and the value of R^2 indicates that the model explains 99.8% of the variability in the reciprocal of the average weight per orange. The



more appropriate model is: $\frac{1}{\text{Ave.} \hat{w}t} = 1.603 + 0.00112(\# \text{Oranges} / \text{Tree})$.

29. Years to live 2008.

- a) The value of R^2 is very high, but the residuals plot shows a clear pattern. This is not a good model for predicting the additional years of life for Hispanic women.
- b) The square root regression model is $\widehat{\sqrt{YrsToLive}} = 9.6867 - 0.0797 \text{ Age}$.

$$\widehat{\sqrt{YrsToLive}} = 9.6867 - 0.0797(18)$$

$$\widehat{\sqrt{YrsToLive}} = 8.2521$$

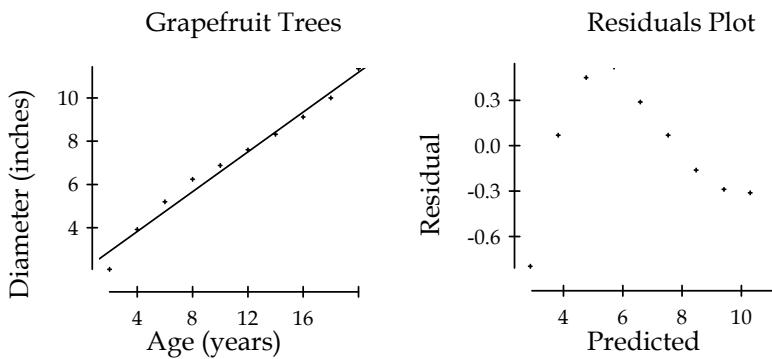
$$\widehat{YrsToLive} = (8.2521)^2 = 68.1 \text{ years}$$

According to the model, the lifespan of an 18-year-old Hispanic women is predicted to be $18 + 68 = 86$ years.

- c) I would have little faith in a prediction made from this model. First of all, the residuals plot showed a pattern, indicating a poor model. Furthermore, the association is about average life expectancy. My friend might have other variables besides age that influence her life expectancy, like being a smoker, having a family history of longevity, and so on.

30. Tree growth.

- a) The association between age and average diameter of grapefruit trees is strong, curved, and positive. Generally, older trees have larger average diameters.



The linear model for this association, $\text{AverageDiameter} = 1.973 + 0.463(\text{Age})$ is not appropriate. The residuals plot shows a clear pattern.

Because of the change in curvature in the association, these data cannot be straightened by re-expression.

- b) If diameters from individual trees were given, instead of averages, the association would have been weaker. Individual observations are more variable than averages.