Regression is one of more widely used statistical methods. It is also widely abused and misinterpreted.

## 8.1 Examining Residuals

**Make sure association is linear before finding the Correlation and the Linear Regression Line.** Sometimes it is hard to tell that a relationship is nonlinear by the Scatterplot, but the Residual plot will show the "curve" better.

Duration of Dive and Dive Heart Rate for Emperor Penquins. $(R^2 = 71.5\%)$
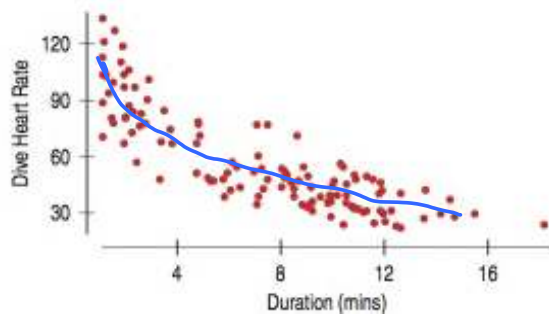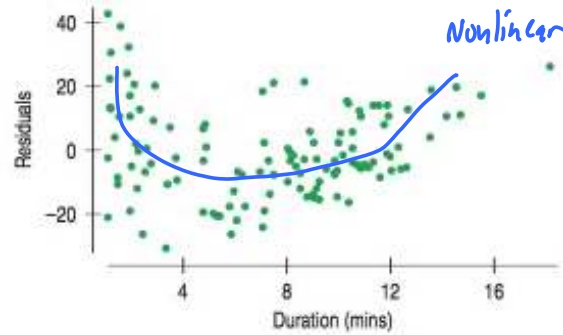
Fig 8.1, p. 215        Fig 8.2, p. 215

### Sifting Residuals for Groups
In Chapter 7 the book looked at an example where the relationship between sugar content and calories of cereals were compared.
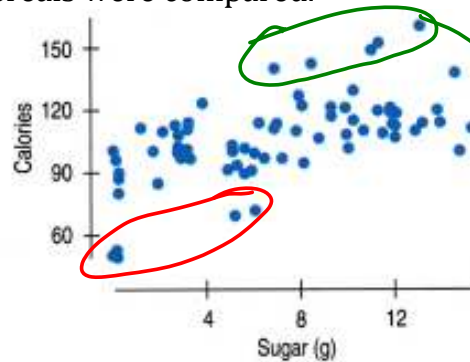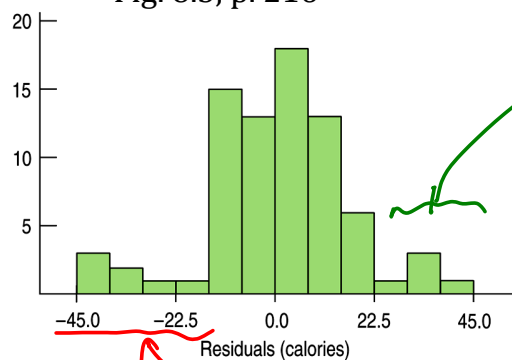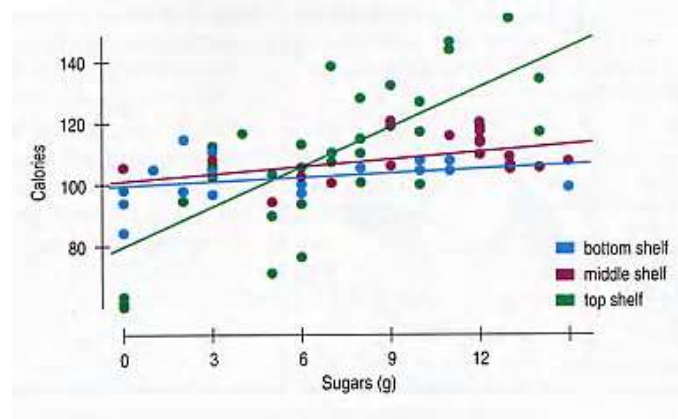
Fig. 8.3, p. 216       Fig. 8.4, p. 216

**Watch out for Subsets in the data.** The data must come from the same group.

Fig 8.5, p. 217



Cereals tend to placed at the eye level of the proposed consumer. Cereals for kids tend to be on the lower shelves while cereals for adults tend to be on the higher shelves.
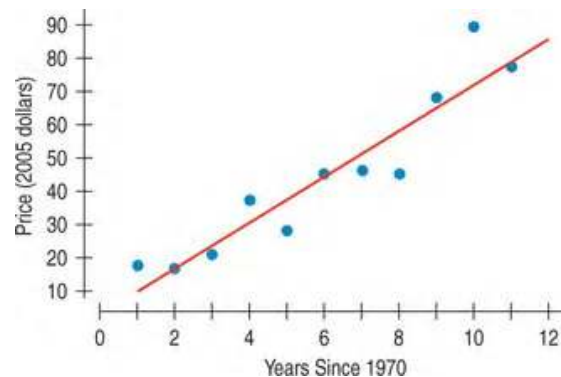
Do you report the regression separately for cereals on the top shelf and cereals on the bottom two shelves?

## 8.2 Extrapolation: Reaching Beyond the Data

**Watch out for Extrapolation.** Extrapolation is when the model is used to predict for values of *x* not in the range of the data. If *x* variable is time then extrapolation is an attempt to see into the future.
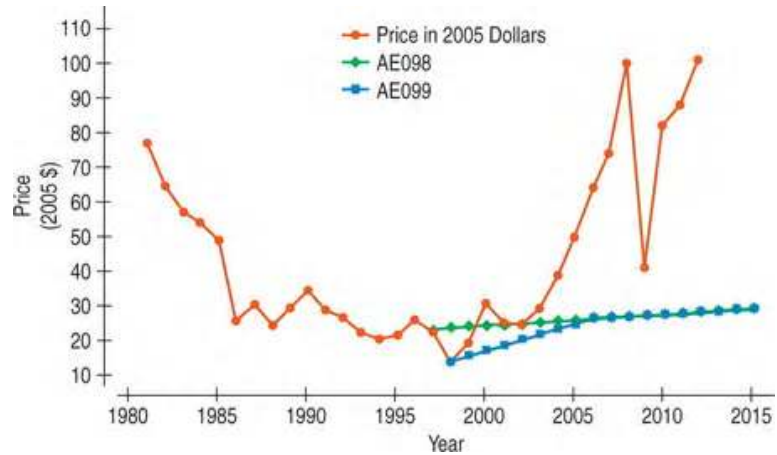
Figure 8.6, p. 218

$$\widehat{Price} = 3.08 + 6.90 \, Years \, Since \, 1970$$



Model from 1971 to 1982 has prices increasing at a rate of $7 per year. Model from 1981 to 1998 shows prices decreased at a rate of $3 per year. Then after 1998 to 2006 there was a sharp increase.

A timeplot from 1981 to 2006 of actual prices along with the Energy Information Administration (EIA) predictions is below.
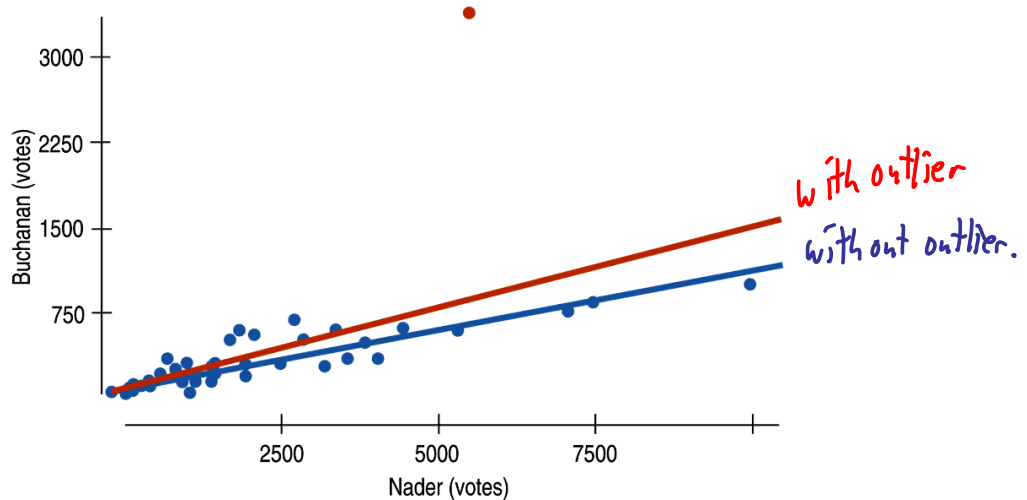
Figure 8.8, p. 219



If extrapolation must be used do not believe it will be correct.

## 8.3 Outliers, Leverage, and Influence

**Watch out for Outliers.**
The following scatterplot shows the number of Nader votes vs number of Buchanan votes for each Florida County in the 2000 Presidential Election (Bush vs Gore).
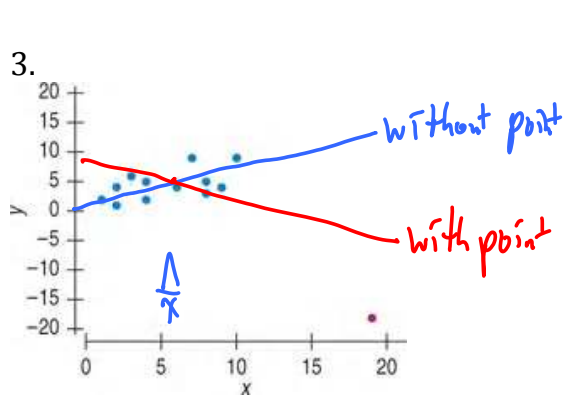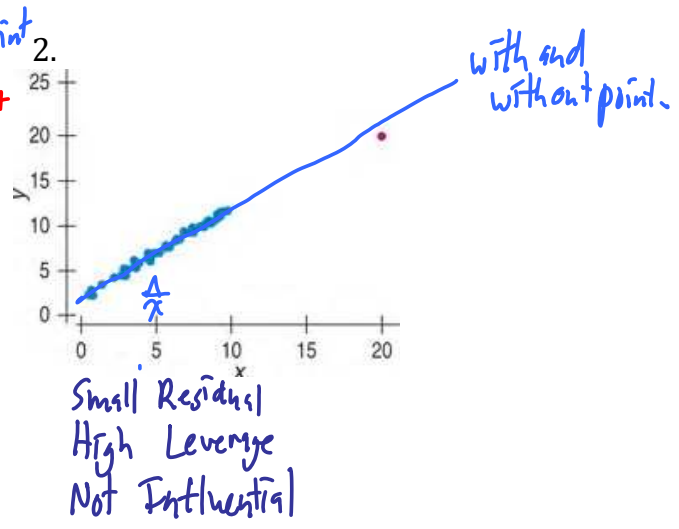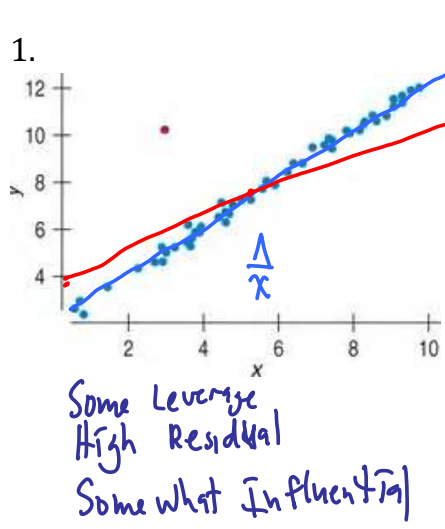
Fig 8.10, p. 222



Red line has slope 0.14 and $R^2 = 42.8\%$. Blue line has slope 0.1 and $R^2 = 82.1\%$

An **Outlier** is any point that stands away from the others.

- Points with **Large Residuals** are far from the regression line. *Vertical Distance*

- Points with **High Leverage** pull the regression line towards them altering the slope and intercept. These points have x-values that are far from $\bar{x}$, think of $\bar{x}$ as being the balancing point for a lever. Their residuals may appear to be small. *Horizontal Distance*
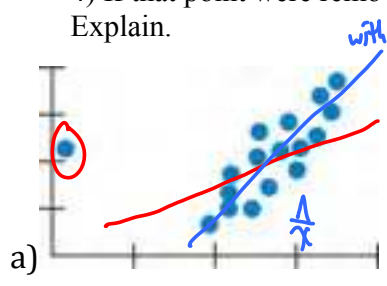
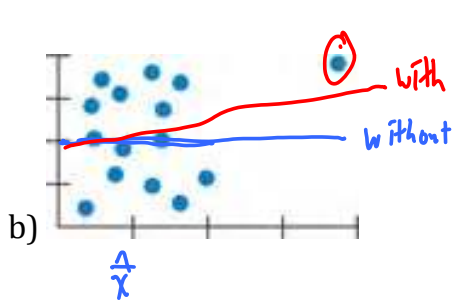A point is considered to be **influential** if omitting the point gives a very different model.

1.



*without point*
*with point*

Some Leverage
High Residual
Some what Influential

2.



*with and without point.*

Small Residual
High Leverage
Not Influential

3.



*without point*
*with point*

Small residual since line pulled toward point.
High Leverage
Influential Point

p.236: **26. More unusual points** Each of the following scatterplots shows a cluster of points and one "stray" point. For each, answer these questions:
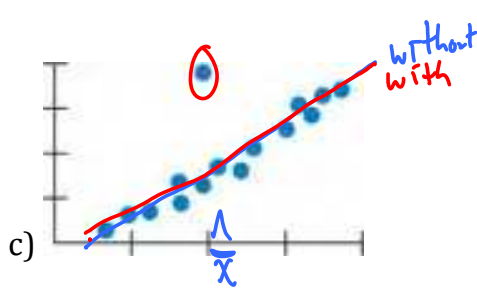
1) In what way is the point unusual? Does it have high leverage, a large residual, or both?
2) Do you think that point is an influential point?
3) If that point were removed, would the correlation become stronger or weaker? Explain.
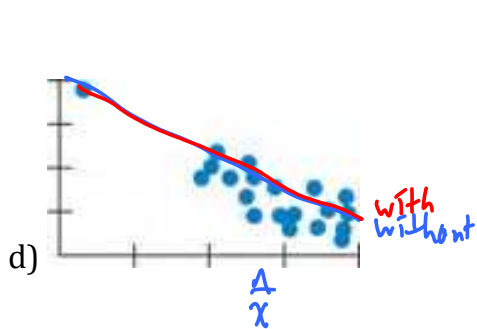4) If that point were removed, would the slope of the regression line increase or decrease? Explain.



*without*

1) Large Residual & High Leverage
2) Influential Point
3) Correlation will be stronger
4) Slope will be increasing.
   will be steeper and more positive

a)

b)



1) Small Residual & High Leverage
2) Influential Point
3) Correlation will be weaker
4) Slope decreases and becomes more horizontal.

c)



1) Large Residual + Small Leverage
2) Not very influential
3) Correlation will be Stronger
4) Slope will be about the Same.

d)



1) Small Residual + High Leverage
2) Not Influential
3) Slightly weaker
4) Slope will be about the same.

## 8.4 Lurking Variables and Causation

**Beware of Lurking Variables.** Lurking Variable is a variable that affects the way the variables in the model appear to be related. No matter how strong the correlation between two variable is this does not give causation (x causes y), because lurking variables can never be ruled out.
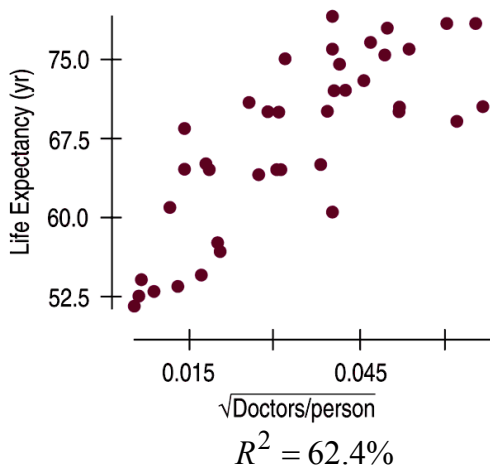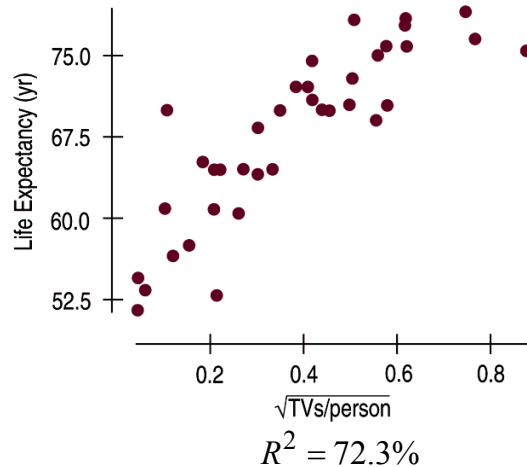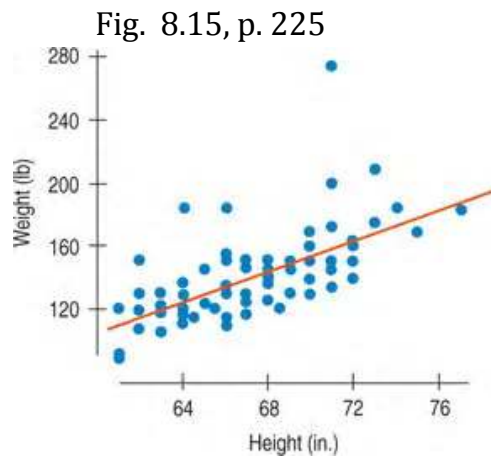
Fig 8.13, p. 224



$R^2 = 62.4\%$

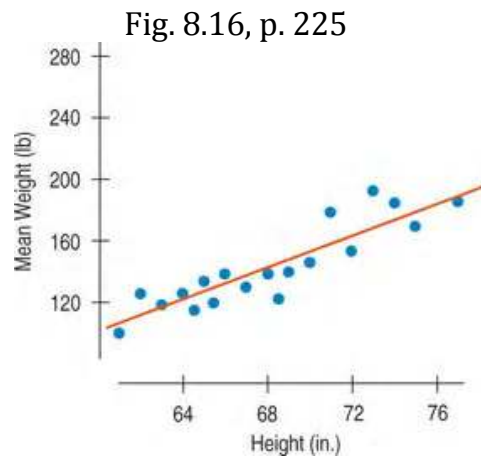Fig 8.14, p. 225



$R^2 = 72.3\%$

Lurking Variable might be income.

## 8.5 Working with Summary Values

**Watch out for data that involves summary values.** Summary values like mean and median will not have the variability as actual data values from individuals.

Fig. 8.15, p. 225

Fig. 8.16, p. 225



$$R^2 = 41.5\%$$

$$R^2 = 80.1\%$$