



Problem 1, continued...

**e.** A 95% confidence interval for  $\mu_Y$  when  $x^* = 70$  is 7702 to 9453. Interpret this interval in words, in the context of this situation.

**f.** The results indicate that counties with higher percentages of high-school graduates tend to have lower crime rates. Can we conclude from this study that having a high school diploma causes people to be less likely to commit crimes, in other words, that higher high-school graduation rates cause crime to be lower? Explain your answer.

**2.** Suppose that a plot of the residuals (Y axis) versus the predicted values (X axis) from a simple linear regression shows a “fan” or “megaphone” shape, with the residuals increasingly spread out as the predicted values increase.

**a.** Which regression assumption would appear to be violated?

**b.** Would it be more appropriate to do a transformation on the X values, or on the Y values? Explain.

**3.** What assumption is being examined by looking at a normal probability plot? Be specific.

**4.** A regression equation is to be fit for predicting  $Y$  = resting pulse rate using the predictor variables  $X_1$  = number of minutes of exercise per week and  $X_2$  = gender, with 1 = male and 0 = female. Explain in words what the coefficient attached to  $X_2$  represents.

**5.** A company offers a training course for the Math SAT. They give their students a test at the end of the course, graded from 0 to 100. They would like to use that test in the future to predict how well students will score on the Math SAT. They have scores on their test and the Math SAT for a sample of students. Thus,  $X$  = score on the company's test and  $Y$  = score on the Math SAT, which ranges from 200 to 800. They plan to use the usual simple linear regression model.

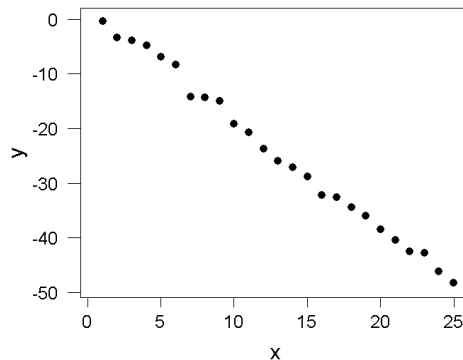
**a.** Would the intercept have a useful meaning in this example? Explain your answer.

**b.** One of the company analysts states that the intercept should be fixed at 200, because that's the lowest the SAT Math score can be. Suppose the intercept is set to 200 for this situation. Write the population model.

**c.** Write the sum that is to be minimized to get the least squares regression line, if the model you wrote in Part b is used.

## MULTIPLE CHOICE

- In a linear regression analysis with the usual assumptions, which one of the following quantities is the same for all individual units in the analysis?
  - $\mu_y$
  - $\sigma_\epsilon$
  - $e_i$
  - $\hat{Y}_i$
- A regression line is used for all of the following *except* one. Which one is *not* a valid use of a regression line in general?
  - to estimate the average value of Y at a specified value of X.
  - to predict the value of Y for an individual, given that individual's X-value.
  - to estimate the change in Y for a one-unit change in X.
  - to determine if a change in X causes a change in Y.
- Which choice is *not* an appropriate description of  $\hat{Y}$  in a regression equation?
  - Estimated response
  - Predicted response
  - Estimated average response
  - Observed response
- Which of the following is the *best* way to determine whether or not there is a statistically significant linear relationship between two quantitative variables?
  - Compute a regression line from a sample and see if the sample slope is 0.
  - Compute the correlation coefficient and see if it is greater than 0.05.
  - Conduct a test of the null hypothesis that the population slope is 0.
  - Conduct a test of the null hypothesis that the population intercept is 0.
- Shown below is a scatterplot of Y versus X.



Which choice is most likely to be the approximate value of  $R^2$ ?

- 99.5%
- 2.0%
- 50.0%
- 99.5%