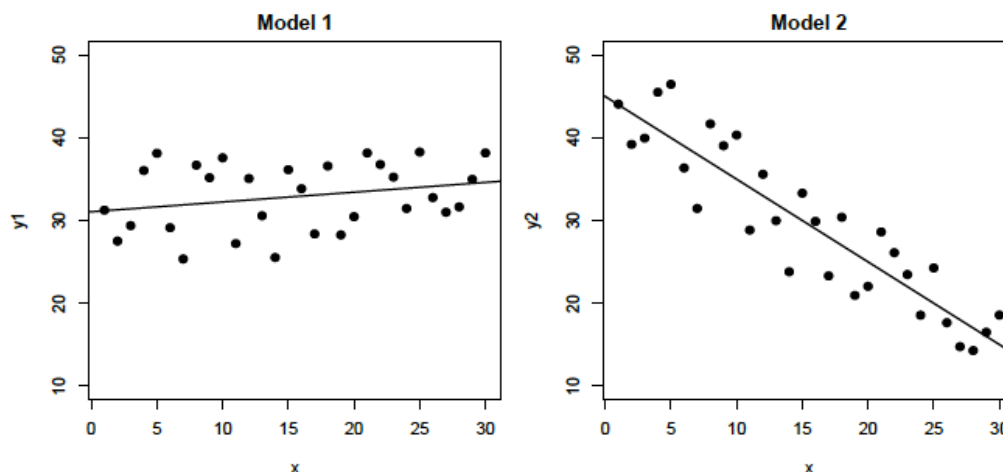PART I: Multiple Choice (5 Points Per Question). Choose the **best** answer.

1. A statistician is comparing two proposed models. The same predictor $x$ is being used for both response variables, $y_1$ and $y_2$. Residual sum of squares is the same for both models. The plots are drawn to scale on the same axes. Which of the following statements is true?
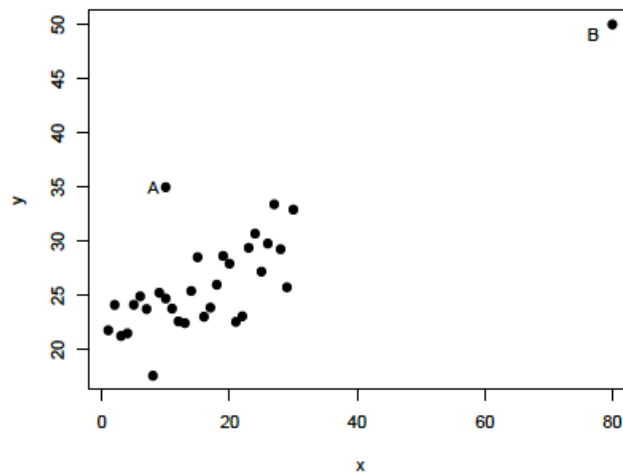


   (a) SSReg for Model 1 is greater than SSReg for Model 2.

   (b) **SSReg for Model 2 is greater than SSReg for Model 1. (NOTE THE DIFFERENCE BETWEEN $\hat{y}_i$ AND $\bar{y}$ IS MUCH LARGER FOR MOST POINTS ON THE SECOND GRAPH.)

   (c) The values of SSReg for both models are equal.

   (d) We know $R^2$ is higher for Model 1, but we can't determine whether SSReg is higher or lower for Model 2.

   (e) We know $R^2$ is higher for Model 2, but we can't determine whether SSReg is higher or lower for Model 1.

2. Which of the following must be true in order for a simple linear regression model to be valid?

   (a) The relationship between $x$ and $y$ must be exponential.

   (b) The relationship between $x$ and $y$ must be quadratic.

   (c) The data must be collected across time.

   (d) **The errors must have constant variance.

   (e) The errors must be normally distributed. (THIS IS NECESSARY ONLY FOR SMALL SAMPLE SIZES AND PREDICTION INTERVALS.)

3. In the Advertisements dataset from class, we looked at the relationship between $x =$ Circulation (in millions) and $y =$ Ad Revenue (in thousands of dollars) of various magazines. The model $\log(\text{AdRevenue}_i) = \beta_0 + \beta_1 \log(\text{Circulation}_i) + e_i$ was fit to the data. The magazine "Family Circle" had a circulation of 3.954, and a confidence interval from the output for that point was $(5.345, 5.458)$. The MSE from the model was 0.0313. How should the confidence interval be back-transformed?

(a) **$\left(e^{5.345+0.0313/2}, e^{5.458+0.0313/2}\right)$ (REMEMBER THAT THE INVERSE TRANSFORMATION FOR LOG IS THE EXPONENTIAL, SO PICK ONE OF THE FIRST TWO AT THE VERY LEAST.)

(b) $\left(e^{5.345-0.0313}, e^{5.458+0.0313}\right)$

(c) $\left(5.345^2 + 0.0313, 5.458^2 + 0.0313\right)$

(d) $\left(5.345^2 - 0.0313, 5.458^2 + 0.0313\right)$

(e) $\left(\frac{1}{5.345}\left(1 + \frac{0.0313}{5.345^2}\right), \frac{1}{5.458}\left(1 + \frac{0.0313}{5.458^2}\right)\right)$

(f) $\left(\frac{1}{5.345}\left(1 - \frac{0.0313}{5.345^2}\right), \frac{1}{5.458}\left(1 + \frac{0.0313}{5.458^2}\right)\right)$

PART II: Multiple Select (2 Points Per Choice). Please circle the letter for **all** of the correct answers; more than one answer may be correct.
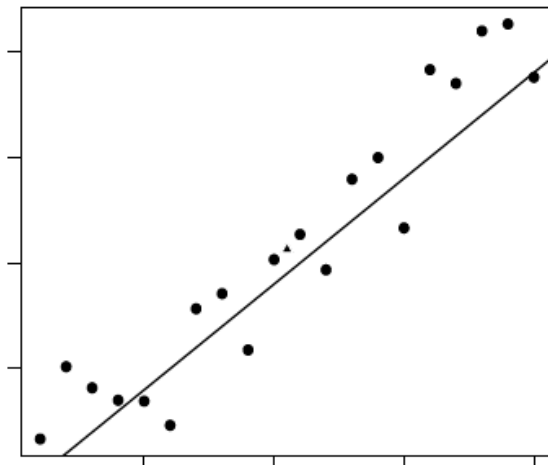
4. The points labeled "A" and "B" in the graph below can be described using which of the following? Select all that apply.



(a) Point A is a high leverage point. (POINT A IS QUITE CLOSE TO THE MEAN OF $x$.)

(b) Point A is influential. (POINT A ISN'T GOING TO CHANGE THE SLOPE SUBSTANTIALLY. IT MIGHT CHANGE THE INTERCEPT A BIT, BUT NOT THE SLOPE.)

(c) **Point B is a high leverage point. (POINT B IS DEFINITELY FAR FROM THE MEAN OF $x$.)

(d) Point B is influential. (POINT B WON'T CHANGE THE SLOPE SUBSTANTIALLY.)

5. Which of the following are reasons for transforming predictor and / or response variables? Select all that apply.

(a) **To ensure the relationship between the predictor and response is a straight line.

(b) **To stabilize the variance of the residuals.

(c) **To reduce the influence of outliers.

(d) **To estimate percentage effects (elasticity).

2

PART III: Short Answer (8 Points Each Part)

6. A model is being fit to some data as shown below; the plot is drawn to scale. The triangle-shaped point in the middle is the point $(\bar{x}, \bar{y})$. Explain to the researcher why the proposed line on the graph is not the least squares regression line.

7. A researcher is interested in the effect of different fertilizers on the amount of corn produced. As a preliminary study, **one** gram of Fertilizer A is added to the first two pots with corn plants, and **two** grams of Fertilizer B are added to the second two pots, for a total of four pots. (Each pot has one corn plant.) The response variable is the weight of the corn produced from the plants. Let $\beta_A$ and $\beta_B$ be the mean amount of corn produced per gram of Fertilizer A and Fertilizer B, respectively.

(a) For the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, where $\mathbf{y}$, $\boldsymbol{\beta}$, and $\mathbf{e}$ are vectors and $\mathbf{X}$ is the design matrix, write out $\mathbf{X}$ and $\boldsymbol{\beta}$.

$$\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 2 \\ 0 & 2 \end{bmatrix}, \; \boldsymbol{\beta} = \begin{bmatrix} \beta_A \\ \beta_B \end{bmatrix}$$

(b) Use the design matrix above and the general least squares solution for the parameter estimate vector to calculate estimates for the parameters.

$$\mathbf{X'X} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 2 & 2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix}$$

3

$$(\mathbf{X'X})^{-1} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.125 \end{bmatrix}$$

$$\mathbf{X'y} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 2 & 2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} y_1 + y_2 \\ 2y_3 + 2y_4 \end{bmatrix}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X'X})^{-1}\mathbf{X'y} = \begin{bmatrix} 0.5\,(y_1 + y_2) \\ 0.25\,(y_3 + y_4) \end{bmatrix}$$

NOTICE THIS MAKES SENSE, SINCE WE'RE INTERESTED IN MEAN GROWTH PER GRAM. THE ESTIMATE FOR MEAN GROWTH PER GRAM FOR FERTILIZER A CAN SIMPLY AVERAGE THE FIRST TWO POTS' GROWTHS, BUT SINCE TWO GRAMS OF FERTILIZER WERE ADDED TO POTS 3 AND 4, WE NEED TO CUT THEIR AVERAGE IN HALF TO GET GROWTH PER GRAM.

PART IV: Long Answer (8 Points Each Part)

8. Education researchers are interested in predicting students who are at risk of failing secondary school in Portugese schools. The response variable of interest is the final year exam grade (*exam*, on a scale of 0-20), and possible predictors include attributes about the students' family lives. We consider two possible predictor variables: the students' current health status (*health*) and how often the student goes out with friends (*goout*), both on a scale of 1: very low to 5: very high. Some relevant R output is shown in the Appendix.

   (a) First we consider a model using health status as a predictor variable. The calculated p-value for the slope of the model is 0.224. At a significance level of 0.05, what does that tell us about using health as a predictor for final year exam performance? Assume all assumptions are met. Explain as if to someone with no statistical experience.

   THE P-VALUE TELLS US THAT WE ARE QUITE LIKELY TO SEE A SLOPE LIKE THE ONE IN OUR DATASET (OR ONE MORE EXTREME) JUST FROM RANDOM CHANCE. WE DON'T HAVE EVIDENCE THAT HEALTH AND EXAM SCORES ARE ASSOCIATED WITH THE POPULATION (BE CAREFUL TO FAIL TO REJECT $H_0 : \beta_1 = 0$ IN YOUR WORDING; DON'T ACCEPT IT!), SO WE WOULDN'T USE IT IF WE WERE AIMING FOR SIMPLICITY.

   (b) Next we consider the model $exam_i = \beta_0 + \beta_1 goout_i + e_i$. Interpret the estimated slope in context, paying special attention to making sure a layman understands what the sign of the slope means.

   WE OBSERVE THAT OUR ESTIMATE OF THE SLOPE IS $\hat{\beta}_1 = -0.5465$. THIS MEANS THAT IF THE FREQUENCY OF STUDENTS' GOING OUT WITH FRIENDS INCREASES BY 1 UNIT (PERHAPS 1 TIME PER WEEK?), WE ESTIMATE THAT THEIR AVERAGE EXAM SCORE DECLINES BY ABOUT HALF A POINT. THE NEGATIVE SLOPE INDICATES THE DECLINE IN SCORES AS GOING OUT INCREASES.

Portugese Schools Output (Going Out Model):

```
Call:
lm(formula = exam ~ go.out)

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  12.1141     0.6793  17.833  < 2e-16
go.out       -0.5465     0.2057  -2.656  0.00823

Residual standard error: 4.547 on 393 degrees of freedom
Multiple R-squared:  0.01763,Adjusted R-squared:  0.01513
F-statistic: 7.054 on 1 and 393 DF,  p-value: 0.008229
```